

Bayesian Linear Regression (Hyperparameter Estimation, Sparse Priors), Bayesian Logistic Regression

Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

Jan 21, 2019



Recap: Bayesian Linear Regression

- Assume Gaussian likelihood: $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \beta^{-1} \mathbf{I}_N)$
- Assume zero-mean spherical Gaussian prior: $p(\mathbf{w}|\lambda) = \prod_{d=1}^D \mathcal{N}(w_d | 0, \lambda^{-1}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda^{-1} \mathbf{I}_D)$



Recap: Bayesian Linear Regression

- Assume Gaussian likelihood: $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \beta^{-1} \mathbf{I}_N)$
- Assume zero-mean spherical Gaussian prior: $p(\mathbf{w}|\lambda) = \prod_{d=1}^D \mathcal{N}(w_d | 0, \lambda^{-1}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda^{-1} \mathbf{I}_D)$
- Assuming hyperparameters as fixed, the posterior is Gaussian

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$



Recap: Bayesian Linear Regression

- Assume Gaussian likelihood: $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \beta^{-1} \mathbf{I}_N)$
- Assume zero-mean spherical Gaussian prior: $p(\mathbf{w}|\lambda) = \prod_{d=1}^D \mathcal{N}(w_d | 0, \lambda^{-1}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda^{-1} \mathbf{I}_D)$
- Assuming hyperparameters as fixed, the posterior is Gaussian

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$

$$\boldsymbol{\Sigma}_N = \left(\beta \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I}_D \right)^{-1} = (\beta \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \quad (\text{posterior's covariance matrix})$$



Recap: Bayesian Linear Regression

- Assume Gaussian likelihood: $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(y_n|\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$
- Assume zero-mean spherical Gaussian prior: $p(\mathbf{w}|\lambda) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$
- Assuming hyperparameters as fixed, the posterior is Gaussian

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$

$$\boldsymbol{\Sigma}_N = \left(\beta \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I}_D \right)^{-1} = \left(\beta \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D \right)^{-1} \quad (\text{posterior's covariance matrix})$$

$$\boldsymbol{\mu}_N = \boldsymbol{\Sigma}_N \left[\beta \sum_{n=1}^N y_n \mathbf{x}_n \right] = \boldsymbol{\Sigma}_N \left[\beta \mathbf{X}^\top \mathbf{y} \right] = \left(\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D \right)^{-1} \mathbf{X}^\top \mathbf{y} \quad (\text{posterior's mean})$$



Recap: Bayesian Linear Regression

- Assume Gaussian likelihood: $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(y_n|\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$
- Assume zero-mean spherical Gaussian prior: $p(\mathbf{w}|\lambda) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$
- Assuming hyperparameters as fixed, the posterior is Gaussian

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$

$$\boldsymbol{\Sigma}_N = \left(\beta \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I}_D \right)^{-1} = \left(\beta \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D \right)^{-1} \quad (\text{posterior's covariance matrix})$$

$$\boldsymbol{\mu}_N = \boldsymbol{\Sigma}_N \left[\beta \sum_{n=1}^N y_n \mathbf{x}_n \right] = \boldsymbol{\Sigma}_N \left[\beta \mathbf{X}^\top \mathbf{y} \right] = \left(\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D \right)^{-1} \mathbf{X}^\top \mathbf{y} \quad (\text{posterior's mean})$$

- The posterior predictive distribution is also Gaussian

$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}, \beta, \lambda) = \int p(y_*|\mathbf{w}, \mathbf{x}_*, \beta) p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) d\mathbf{w} = \mathcal{N}(\boldsymbol{\mu}_N^\top \mathbf{x}_*, \beta^{-1} + \mathbf{x}_*^\top \boldsymbol{\Sigma}_N \mathbf{x}_*)$$



Recap: Bayesian Linear Regression

- Assume Gaussian likelihood: $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(y_n|\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$
- Assume zero-mean spherical Gaussian prior: $p(\mathbf{w}|\lambda) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$
- Assuming hyperparameters as fixed, the posterior is Gaussian

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$

$$\boldsymbol{\Sigma}_N = \left(\beta \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I}_D \right)^{-1} = \left(\beta \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D \right)^{-1} \quad (\text{posterior's covariance matrix})$$

$$\boldsymbol{\mu}_N = \boldsymbol{\Sigma}_N \left[\beta \sum_{n=1}^N y_n \mathbf{x}_n \right] = \boldsymbol{\Sigma}_N \left[\beta \mathbf{X}^\top \mathbf{y} \right] = \left(\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D \right)^{-1} \mathbf{X}^\top \mathbf{y} \quad (\text{posterior's mean})$$

- The posterior predictive distribution is also Gaussian

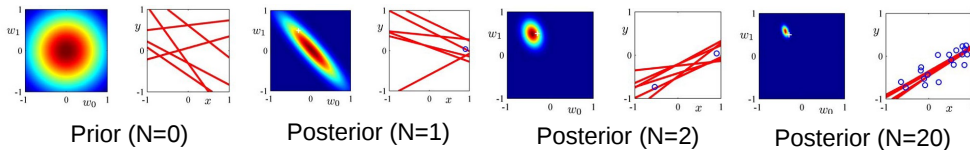
$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}, \beta, \lambda) = \int p(y_*|\mathbf{w}, \mathbf{x}_*, \beta) p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) d\mathbf{w} = \mathcal{N}(\boldsymbol{\mu}_N^\top \mathbf{x}_*, \beta^{-1} + \mathbf{x}_*^\top \boldsymbol{\Sigma}_N \mathbf{x}_*)$$

- Gives both **predictive mean** and **predictive variance** (imp: pred-var is different for each input)



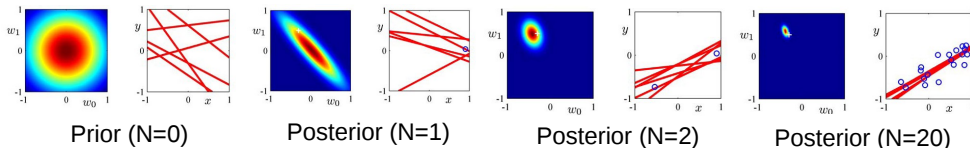
A Visualization of Uncertainty in Bayesian Linear Regression

- Posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ and lines (w_0 intercept, w_1 slope) corresponding to some random \mathbf{w} 's

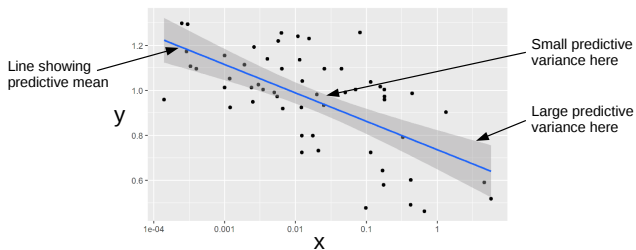


A Visualization of Uncertainty in Bayesian Linear Regression

- Posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ and lines (w_0 intercept, w_1 slope) corresponding to some random \mathbf{w} 's

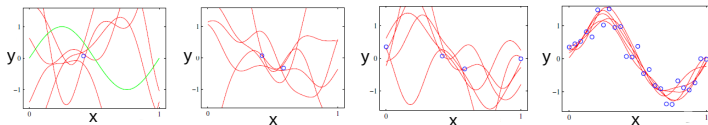


- A visualization of the posterior predictive of a Bayesian linear regression model



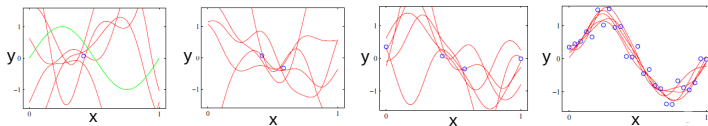
A Visualization of Uncertainty (Contd)

- We can similarly visualize a Bayesian nonlinear regression model
- Figures below: Green curve is the true function and blue circles are observations (x_n, y_n)
- Posterior of the nonlinear regression model: Some curves drawn from the posterior

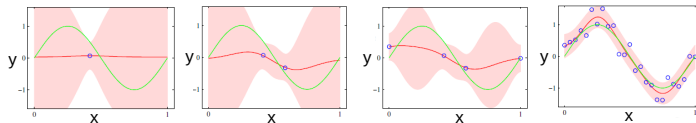


A Visualization of Uncertainty (Contd)

- We can similarly visualize a Bayesian nonlinear regression model
- Figures below: Green curve is the true function and blue circles are observations (x_n, y_n)
- Posterior of the nonlinear regression model: Some curves drawn from the posterior



- Posterior predictive: Red curve is predictive mean, shaded region denotes predictive uncertainty



Estimating Hyperparameters for Bayesian Linear Regression



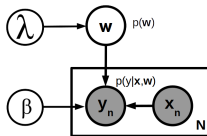
Learning Hyperparameters in Probabilistic Models

- Can treat hyperparams as just **a bunch of additional unknowns**
- Can be learned using a suitable inference algorithm (point estimation or fully Bayesian)



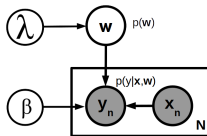
Learning Hyperparameters in Probabilistic Models

- Can treat hyperparams as just a bunch of additional unknowns
- Can be learned using a suitable inference algorithm (point estimation or fully Bayesian)
- Example: For the linear regression model, the full set of parameters would be $(\mathbf{w}, \lambda, \beta)$



Learning Hyperparameters in Probabilistic Models

- Can treat hyperparams as just **a bunch of additional unknowns**
- Can be learned using a suitable inference algorithm (point estimation or fully Bayesian)
- Example: For the linear regression model, the full set of parameters would be $(\mathbf{w}, \lambda, \beta)$



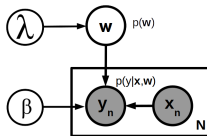
- Can assume priors on all these parameters and infer their “joint” posterior distribution

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w}, \lambda, \beta)}{p(\mathbf{y} | \mathbf{X})} = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w} | \lambda) p(\beta) p(\lambda)}{\int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \lambda) p(\beta) p(\lambda) d\mathbf{w} d\lambda d\beta}$$



Learning Hyperparameters in Probabilistic Models

- Can treat hyperparams as just **a bunch of additional unknowns**
- Can be learned using a suitable inference algorithm (point estimation or fully Bayesian)
- Example: For the linear regression model, the full set of parameters would be $(\mathbf{w}, \lambda, \beta)$



- Can assume priors on all these parameters and infer their “joint” posterior distribution

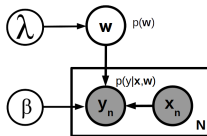
$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w}, \lambda, \beta)}{p(\mathbf{y} | \mathbf{X})} = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w} | \lambda) p(\beta) p(\lambda)}{\int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \lambda) p(\beta) p(\lambda) d\mathbf{w} d\lambda d\beta}$$

- Inferring the above is **usually intractable** (rare to have conjugacy).



Learning Hyperparameters in Probabilistic Models

- Can treat hyperparams as just **a bunch of additional unknowns**
- Can be learned using a suitable inference algorithm (point estimation or fully Bayesian)
- Example: For the linear regression model, the full set of parameters would be $(\mathbf{w}, \lambda, \beta)$



- Can assume priors on all these parameters and infer their “joint” posterior distribution

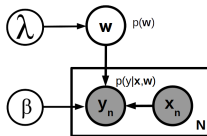
$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w}, \lambda, \beta)}{p(\mathbf{y} | \mathbf{X})} = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w} | \lambda) p(\beta) p(\lambda)}{\int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \lambda) p(\beta) p(\lambda) d\mathbf{w} d\lambda d\beta}$$

- Inferring the above is **usually intractable** (rare to have conjugacy). Requires approximations.



Learning Hyperparameters in Probabilistic Models

- Can treat hyperparams as just **a bunch of additional unknowns**
- Can be learned using a suitable inference algorithm (point estimation or fully Bayesian)
- Example: For the linear regression model, the full set of parameters would be $(\mathbf{w}, \lambda, \beta)$



- Can assume priors on all these parameters and infer their “joint” posterior distribution

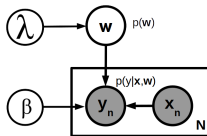
$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w}, \lambda, \beta)}{p(\mathbf{y} | \mathbf{X})} = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w} | \lambda) p(\beta) p(\lambda)}{\int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \lambda) p(\beta) p(\lambda) d\mathbf{w} d\lambda d\beta}$$

- Inferring the above is **usually intractable** (rare to have conjugacy). Requires approximations. Also,



Learning Hyperparameters in Probabilistic Models

- Can treat hyperparams as just **a bunch of additional unknowns**
- Can be learned using a suitable inference algorithm (point estimation or fully Bayesian)
- Example: For the linear regression model, the full set of parameters would be $(\mathbf{w}, \lambda, \beta)$



- Can assume priors on all these parameters and infer their “joint” posterior distribution

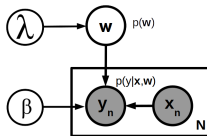
$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w}, \lambda, \beta)}{p(\mathbf{y} | \mathbf{X})} = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w} | \lambda) p(\beta) p(\lambda)}{\int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \lambda) p(\beta) p(\lambda) d\mathbf{w} d\lambda d\beta}$$

- Inferring the above is **usually intractable** (rare to have conjugacy). Requires approximations. Also,
 - What priors (or “hyperpriors”) to choose for β and λ ?



Learning Hyperparameters in Probabilistic Models

- Can treat hyperparams as just **a bunch of additional unknowns**
- Can be learned using a suitable inference algorithm (point estimation or fully Bayesian)
- Example: For the linear regression model, the full set of parameters would be $(\mathbf{w}, \lambda, \beta)$



- Can assume priors on all these parameters and infer their “joint” posterior distribution

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w}, \lambda, \beta)}{p(\mathbf{y} | \mathbf{X})} = \frac{p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta, \lambda) p(\mathbf{w} | \lambda) p(\beta) p(\lambda)}{\int p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} | \lambda) p(\beta) p(\lambda) d\mathbf{w} d\lambda d\beta}$$

- Inferring the above is **usually intractable** (rare to have conjugacy). Requires approximations. Also,
 - What priors (or “hyperpriors”) to choose for β and λ ?
 - What about the hyperparameters of those priors?



Learning Hyperparameters via Point Estimation

- One popular way to estimate hyperparameters is by maximizing the **marginal likelihood**



Learning Hyperparameters via Point Estimation

- One popular way to estimate hyperparameters is by maximizing the **marginal likelihood**
- For our linear regression model, this quantity (a function of the hyperparams) will be

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)p(\mathbf{w}|\lambda)d\mathbf{w}$$



Learning Hyperparameters via Point Estimation

- One popular way to estimate hyperparameters is by maximizing the **marginal likelihood**
- For our linear regression model, this quantity (a function of the hyperparams) will be

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- The “optimal” hyperparameters in this case can be then found by

$$\hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} \log p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$$



Learning Hyperparameters via Point Estimation

- One popular way to estimate hyperparameters is by maximizing the **marginal likelihood**
- For our linear regression model, this quantity (a function of the hyperparams) will be

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- The “optimal” hyperparameters in this case can be then found by

$$\hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} \log p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$$

- This is called **MLE-II** or (log) evidence maximization



Learning Hyperparameters via Point Estimation

- One popular way to estimate hyperparameters is by maximizing the **marginal likelihood**
- For our linear regression model, this quantity (a function of the hyperparams) will be

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- The “optimal” hyperparameters in this case can be then found by

$$\hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} \log p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$$

- This is called **MLE-II** or (log) evidence maximization
 - Akin to doing MLE to estimate the hyperparameters where the “main” parameter (in this case \mathbf{w}) has been integrated out from the model's likelihood function



Learning Hyperparameters via Point Estimation

- One popular way to estimate hyperparameters is by maximizing the **marginal likelihood**
- For our linear regression model, this quantity (a function of the hyperparams) will be

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- The “optimal” hyperparameters in this case can be then found by

$$\hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} \log p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$$

- This is called **MLE-II** or (log) evidence maximization
 - Akin to doing MLE to estimate the hyperparameters where the “main” parameter (in this case \mathbf{w}) has been integrated out from the model's likelihood function
- **Note:** If the likelihood and prior are conjugate then marginal likelihood is available in closed form



What is MLE-II Doing?

- For linear regression case, would ideally like the posterior over all unknowns, i.e., $p(\mathbf{w}, \lambda, \beta | \mathbf{X}, \mathbf{y})$

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \quad (\text{from product rule})$$



What is MLE-II Doing?

- For linear regression case, would ideally like the posterior over all unknowns, i.e., $p(\mathbf{w}, \lambda, \beta | \mathbf{X}, \mathbf{y})$

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \quad (\text{from product rule})$$

- Note that $p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda)$ is easy if λ, β are known



What is MLE-II Doing?

- For linear regression case, would ideally like the posterior over all unknowns, i.e., $p(\mathbf{w}, \lambda, \beta | \mathbf{X}, \mathbf{y})$

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \quad (\text{from product rule})$$

- Note that $p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda)$ is easy if λ, β are known
- However $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\beta) p(\lambda)}{p(\mathbf{y} | \mathbf{X})}$ is hard



What is MLE-II Doing?

- For linear regression case, would ideally like the posterior over all unknowns, i.e., $p(\mathbf{w}, \lambda, \beta | \mathbf{X}, \mathbf{y})$

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \quad (\text{from product rule})$$

- Note that $p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda)$ is easy if λ, β are known
- However $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\beta) p(\lambda)}{p(\mathbf{y} | \mathbf{X})}$ is hard (lack of conjugacy, intractable denominator)



What is MLE-II Doing?

- For linear regression case, would ideally like the posterior over all unknowns, i.e., $p(\mathbf{w}, \lambda, \beta | \mathbf{X}, \mathbf{y})$

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \quad (\text{from product rule})$$

- Note that $p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda)$ is easy if λ, β are known
- However $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\beta) p(\lambda)}{p(\mathbf{y} | \mathbf{X})}$ is hard (lack of conjugacy, intractable denominator)
- Let's **approximate** it by a **point function** δ at the **mode** of $p(\beta, \lambda | \mathbf{X}, \mathbf{y})$

$$p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda})$$



What is MLE-II Doing?

- For linear regression case, would ideally like the posterior over all unknowns, i.e., $p(\mathbf{w}, \lambda, \beta | \mathbf{X}, \mathbf{y})$

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \quad (\text{from product rule})$$

- Note that $p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda)$ is easy if λ, β are known
- However $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\beta) p(\lambda)}{p(\mathbf{y} | \mathbf{X})}$ is hard (lack of conjugacy, intractable denominator)
- Let's **approximate** it by a **point function** δ at the **mode** of $p(\beta, \lambda | \mathbf{X}, \mathbf{y})$

$$p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda}) \quad \text{where} \quad \hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} p(\beta, \lambda | \mathbf{X}, \mathbf{y}) = \arg \max_{\beta, \lambda} p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\lambda) p(\beta)$$



What is MLE-II Doing?

- For linear regression case, would ideally like the posterior over all unknowns, i.e., $p(\mathbf{w}, \lambda, \beta | \mathbf{X}, \mathbf{y})$

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \quad (\text{from product rule})$$

- Note that $p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda)$ is easy if λ, β are known
- However $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\beta) p(\lambda)}{p(\mathbf{y} | \mathbf{X})}$ is hard (lack of conjugacy, intractable denominator)
- Let's **approximate** it by a **point function** δ at the **mode** of $p(\beta, \lambda | \mathbf{X}, \mathbf{y})$

$$p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda}) \quad \text{where} \quad \hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} p(\beta, \lambda | \mathbf{X}, \mathbf{y}) = \arg \max_{\beta, \lambda} p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\lambda) p(\beta)$$

- Moreover, if $p(\beta), p(\lambda)$ are uniform/uninformative priors then

$$\hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} p(\mathbf{y} | \mathbf{X}, \beta, \lambda)$$



What is MLE-II Doing?

- For linear regression case, would ideally like the posterior over all unknowns, i.e., $p(\mathbf{w}, \lambda, \beta | \mathbf{X}, \mathbf{y})$

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \quad (\text{from product rule})$$

- Note that $p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda)$ is easy if λ, β are known
- However $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\beta) p(\lambda)}{p(\mathbf{y} | \mathbf{X})}$ is hard (lack of conjugacy, intractable denominator)
- Let's **approximate** it by a **point function** δ at the **mode** of $p(\beta, \lambda | \mathbf{X}, \mathbf{y})$

$$p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda}) \quad \text{where} \quad \hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} p(\beta, \lambda | \mathbf{X}, \mathbf{y}) = \arg \max_{\beta, \lambda} p(\mathbf{y} | \mathbf{X}, \beta, \lambda) p(\lambda) p(\beta)$$

- Moreover, if $p(\beta), p(\lambda)$ are uniform/uninformative priors then

$$\hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} p(\mathbf{y} | \mathbf{X}, \beta, \lambda)$$

- Thus MLE-II is approximating the posterior of hyperparams by their point estimate assuming uniform priors (therefore we don't need to worry about a prior over the hyperparams)



MLE-II for Linear Regression

- For the linear regression case, the marginal likelihood is defined as

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$



MLE-II for Linear Regression

- For the linear regression case, the marginal likelihood is defined as

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- Since $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$ and $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$, the marginal likelihood

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top)$$



MLE-II for Linear Regression

- For the linear regression case, the marginal likelihood is defined as

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- Since $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$ and $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$, the marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \beta, \lambda) &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top) \\ &= \frac{1}{(2\pi)^{N/2}} |\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{y}^\top (\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y}\right) \end{aligned}$$



MLE-II for Linear Regression

- For the linear regression case, the marginal likelihood is defined as

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- Since $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$ and $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$, the marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \beta, \lambda) &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top) \\ &= \frac{1}{(2\pi)^{N/2}} |\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{y}^\top (\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y}\right) \end{aligned}$$

- MLE-II maximizes $\log p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$ w.r.t. β and λ to estimate these hyperparams



MLE-II for Linear Regression

- For the linear regression case, the marginal likelihood is defined as

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- Since $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$ and $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$, the marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \beta, \lambda) &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top) \\ &= \frac{1}{(2\pi)^{N/2}} |\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{y}^\top (\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y}\right) \end{aligned}$$

- MLE-II maximizes $\log p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$ w.r.t. β and λ to estimate these hyperparams
 - This objective doesn't have a closed form solution



MLE-II for Linear Regression

- For the linear regression case, the marginal likelihood is defined as

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- Since $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$ and $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$, the marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \beta, \lambda) &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top) \\ &= \frac{1}{(2\pi)^{N/2}} |\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{y}^\top (\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y}\right) \end{aligned}$$

- MLE-II maximizes $\log p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$ w.r.t. β and λ to estimate these hyperparams
 - This objective doesn't have a closed form solution
 - Solved using iterative/alternating optimization



MLE-II for Linear Regression

- For the linear regression case, the marginal likelihood is defined as

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- Since $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$ and $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$, the marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \beta, \lambda) &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top) \\ &= \frac{1}{(2\pi)^{N/2}} |\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{y}^\top (\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y}\right) \end{aligned}$$

- MLE-II maximizes $\log p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$ w.r.t. β and λ to estimate these hyperparams
 - This objective doesn't have a closed form solution
 - Solved using iterative/alternating optimization
 - PRML Chapter 3 contains the iterative update equations



MLE-II for Linear Regression

- For the linear regression case, the marginal likelihood is defined as

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- Since $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$ and $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$, the marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \beta, \lambda) &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top) \\ &= \frac{1}{(2\pi)^{N/2}} |\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{y}^\top (\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y}\right) \end{aligned}$$

- MLE-II maximizes $\log p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$ w.r.t. β and λ to estimate these hyperparams
 - This objective doesn't have a closed form solution
 - Solved using iterative/alternating optimization
 - PRML Chapter 3 contains the iterative update equations
- Note: Can also do “MAP-II” using a suitable prior on these hyperparams (e.g., gamma)



MLE-II for Linear Regression

- For the linear regression case, the marginal likelihood is defined as

$$p(\mathbf{y}|\mathbf{X}, \beta, \lambda) = \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w}|\lambda) d\mathbf{w}$$

- Since $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$ and $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)$, the marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \beta, \lambda) &= \mathcal{N}(\mathbf{y}|\mathbf{0}, \beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top) \\ &= \frac{1}{(2\pi)^{N/2}} |\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top|^{-1/2} \exp\left(-\frac{1}{2}\mathbf{y}^\top (\beta^{-1}\mathbf{I} + \lambda^{-1}\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{y}\right) \end{aligned}$$

- MLE-II maximizes $\log p(\mathbf{y}|\mathbf{X}, \beta, \lambda)$ w.r.t. β and λ to estimate these hyperparams
 - This objective doesn't have a closed form solution
 - Solved using iterative/alternating optimization
 - PRML Chapter 3 contains the iterative update equations
- Note: Can also do "MAP-II" using a suitable prior on these hyperparams (e.g., gamma)
- Note: Can also use different λ_d for each w_d



Using MLE-II Estimates for Making Prediction

- With the MLE-II approximation $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda})$, the posterior over unknowns

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y})$$



Using MLE-II Estimates for Making Prediction

- With the MLE-II approximation $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda})$, the posterior over unknowns

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda})$$



Using MLE-II Estimates for Making Prediction

- With the MLE-II approximation $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda})$, the posterior over unknowns

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda})$$

- The posterior predictive distribution can also be approximated as

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) d\mathbf{w} d\beta d\lambda$$



Using MLE-II Estimates for Making Prediction

- With the MLE-II approximation $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda})$, the posterior over unknowns

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda})$$

- The posterior predictive distribution can also be approximated as

$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) d\mathbf{w} d\beta d\lambda \\ &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) d\beta d\lambda d\mathbf{w} \end{aligned}$$



Using MLE-II Estimates for Making Prediction

- With the MLE-II approximation $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda})$, the posterior over unknowns

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda})$$

- The posterior predictive distribution can also be approximated as

$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) d\mathbf{w} d\beta d\lambda \\ &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) d\beta d\lambda d\mathbf{w} \\ &\approx \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda}) d\mathbf{w} \end{aligned}$$



Using MLE-II Estimates for Making Prediction

- With the MLE-II approximation $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda})$, the posterior over unknowns

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) = p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda})$$

- The posterior predictive distribution can also be approximated as

$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) d\mathbf{w} d\beta d\lambda \\ &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) d\beta d\lambda d\mathbf{w} \\ &\approx \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda}) d\mathbf{w} \end{aligned}$$

- This is also the same as the usual posterior predictive distribution we have seen earlier, except we are treating the hyperparams $\hat{\beta}, \hat{\lambda}$ fixed at their MLE-II based estimates



Modeling Sparse Weights



Modeling Sparse Weights

- Many probabilistic models consist of weights that are given zero-mean Gaussian priors, e.g.,

$$\mu(\mathbf{x}) = \sum_{d=1}^D w_d x_d \quad (\text{mean of a prob. lin reg model})$$

$$\mu(\mathbf{x}) = \sum_{n=1}^N w_n k(\mathbf{x}_n, \mathbf{x}) \quad (\text{mean of a prob. kernel based nonlin reg model})$$



Modeling Sparse Weights

- Many probabilistic models consist of weights that are given zero-mean Gaussian priors, e.g.,

$$\mu(\mathbf{x}) = \sum_{d=1}^D w_d x_d \quad (\text{mean of a prob. lin reg model})$$

$$\mu(\mathbf{x}) = \sum_{n=1}^N w_n k(\mathbf{x}_n, \mathbf{x}) \quad (\text{mean of a prob. kernel based nonlin reg model})$$

- A zero-mean prior is of the form $p(w_d) = \mathcal{N}(0, \lambda^{-1})$ or $p(w_d) = \mathcal{N}(0, \lambda_d^{-1})$



Modeling Sparse Weights

- Many probabilistic models consist of weights that are given zero-mean Gaussian priors, e.g.,

$$\mu(\mathbf{x}) = \sum_{d=1}^D w_d x_d \quad (\text{mean of a prob. lin reg model})$$

$$\mu(\mathbf{x}) = \sum_{n=1}^N w_n k(\mathbf{x}_n, \mathbf{x}) \quad (\text{mean of a prob. kernel based nonlin reg model})$$

- A zero-mean prior is of the form $p(w_d) = \mathcal{N}(0, \lambda^{-1})$ or $p(w_d) = \mathcal{N}(0, \lambda_d^{-1})$
- Precision λ or λ_d specifies our belief about how close to zero w_d is (like regularization hyperparam)



Modeling Sparse Weights

- Many probabilistic models consist of weights that are given zero-mean Gaussian priors, e.g.,

$$\mu(\mathbf{x}) = \sum_{d=1}^D w_d x_d \quad (\text{mean of a prob. lin reg model})$$

$$\mu(\mathbf{x}) = \sum_{n=1}^N w_n k(\mathbf{x}_n, \mathbf{x}) \quad (\text{mean of a prob. kernel based nonlin reg model})$$

- A zero-mean prior is of the form $p(w_d) = \mathcal{N}(0, \lambda^{-1})$ or $p(w_d) = \mathcal{N}(0, \lambda_d^{-1})$
- Precision λ or λ_d specifies our belief about how close to zero w_d is (like regularization hyperparam)
- However, such a prior usually gives small weights but not very strong sparsity



Modeling Sparse Weights

- Many probabilistic models consist of weights that are given zero-mean Gaussian priors, e.g.,

$$\mu(\mathbf{x}) = \sum_{d=1}^D w_d x_d \quad (\text{mean of a prob. lin reg model})$$

$$\mu(\mathbf{x}) = \sum_{n=1}^N w_n k(\mathbf{x}_n, \mathbf{x}) \quad (\text{mean of a prob. kernel based nonlin reg model})$$

- A zero-mean prior is of the form $p(w_d) = \mathcal{N}(0, \lambda^{-1})$ or $p(w_d) = \mathcal{N}(0, \lambda_d^{-1})$
- Precision λ or λ_d specifies our belief about how close to zero w_d is (like regularization hyperparam)
- However, such a prior usually gives small weights but not very strong sparsity
- Putting a gamma prior on precision can give **sparsity** (will soon see why)



Modeling Sparse Weights

- Many probabilistic models consist of weights that are given zero-mean Gaussian priors, e.g.,

$$\mu(\mathbf{x}) = \sum_{d=1}^D w_d x_d \quad (\text{mean of a prob. lin reg model})$$

$$\mu(\mathbf{x}) = \sum_{n=1}^N w_n k(\mathbf{x}_n, \mathbf{x}) \quad (\text{mean of a prob. kernel based nonlin reg model})$$

- A zero-mean prior is of the form $p(w_d) = \mathcal{N}(0, \lambda^{-1})$ or $p(w_d) = \mathcal{N}(0, \lambda_d^{-1})$
- Precision λ or λ_d specifies our belief about how close to zero w_d is (like regularization hyperparam)
- However, such a prior usually gives small weights but not very strong sparsity
- Putting a gamma prior on precision can give **sparsity** (will soon see why)
- Sparsity of weights will be a very useful thing to have in many models, e.g.,



Modeling Sparse Weights

- Many probabilistic models consist of weights that are given zero-mean Gaussian priors, e.g.,

$$\mu(\mathbf{x}) = \sum_{d=1}^D w_d x_d \quad (\text{mean of a prob. lin reg model})$$

$$\mu(\mathbf{x}) = \sum_{n=1}^N w_n k(\mathbf{x}_n, \mathbf{x}) \quad (\text{mean of a prob. kernel based nonlin reg model})$$

- A zero-mean prior is of the form $p(w_d) = \mathcal{N}(0, \lambda^{-1})$ or $p(w_d) = \mathcal{N}(0, \lambda_d^{-1})$
- Precision λ or λ_d specifies our belief about how close to zero w_d is (like regularization hyperparam)
- However, such a prior usually gives small weights but not very strong sparsity
- Putting a gamma prior on precision can give **sparsity** (will soon see why)
- Sparsity of weights will be a very useful thing to have in many models, e.g.,
 - For linear model, this helps learn relevance of each feature x_d



Modeling Sparse Weights

- Many probabilistic models consist of weights that are given zero-mean Gaussian priors, e.g.,

$$\mu(\mathbf{x}) = \sum_{d=1}^D w_d x_d \quad (\text{mean of a prob. lin reg model})$$

$$\mu(\mathbf{x}) = \sum_{n=1}^N w_n k(\mathbf{x}_n, \mathbf{x}) \quad (\text{mean of a prob. kernel based nonlin reg model})$$

- A zero-mean prior is of the form $p(w_d) = \mathcal{N}(0, \lambda^{-1})$ or $p(w_d) = \mathcal{N}(0, \lambda_d^{-1})$
- Precision λ or λ_d specifies our belief about how close to zero w_d is (like regularization hyperparam)
- However, such a prior usually gives small weights but not very strong sparsity
- Putting a gamma prior on precision can give **sparsity** (will soon see why)
- Sparsity of weights will be a very useful thing to have in many models, e.g.,
 - For linear model, this helps learn relevance of each feature x_d
 - For kernel based model, this helps learn the relevance of each input \mathbf{x}_n (**Relevance Vector Machine**)



Sparsity via a Hierarchical Prior

- Consider linear regression with prior $p(w_d|\lambda_d) = \mathcal{N}(0, \lambda_d^{-1})$ on each weight
- Let's treat precision λ_d as unknown and use a gamma (shape = a , rate = b) prior on it

$$p(\lambda_d) = \text{Gamma}(a, b) = \frac{b^a}{\Gamma(a)} \lambda_d^{a-1} \exp(-b\lambda_d)$$



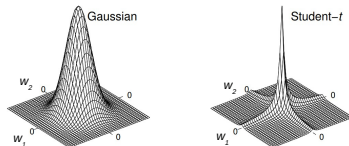
Sparsity via a Hierarchical Prior

- Consider linear regression with prior $p(w_d|\lambda_d) = \mathcal{N}(0, \lambda_d^{-1})$ on each weight
- Let's treat precision λ_d as unknown and use a gamma (shape = a , rate = b) prior on it

$$p(\lambda_d) = \text{Gamma}(a, b) = \frac{b^a}{\Gamma(a)} \lambda_d^{a-1} \exp(-b\lambda_d)$$

- Marginalizing the precision leads to a **Student-t prior** on each w_d

$$p(w_d) = \int p(w_d|\lambda_d)p(\lambda_d)d\lambda_d = \frac{b^a \Gamma(a + 1/2)}{\sqrt{2\pi} \Gamma(a)} (b + w_d^2/2)^{-(a+1/2)}$$



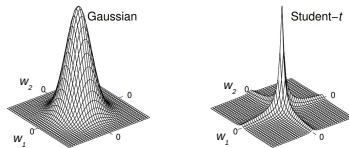
Sparsity via a Hierarchical Prior

- Consider linear regression with prior $p(w_d|\lambda_d) = \mathcal{N}(0, \lambda_d^{-1})$ on each weight
- Let's treat precision λ_d as unknown and use a gamma (shape = a , rate = b) prior on it

$$p(\lambda_d) = \text{Gamma}(a, b) = \frac{b^a}{\Gamma(a)} \lambda_d^{a-1} \exp(-b\lambda_d)$$

- Marginalizing the precision leads to a **Student-t prior** on each w_d

$$p(w_d) = \int p(w_d|\lambda_d)p(\lambda_d)d\lambda_d = \frac{b^a \Gamma(a + 1/2)}{\sqrt{2\pi} \Gamma(a)} (b + w_d^2/2)^{-(a+1/2)}$$



- Note: Can make the prior an **uninformative prior** by setting a and b to be very small (e.g., 10^{-4})

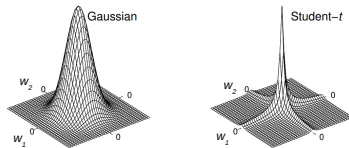
Sparsity via a Hierarchical Prior

- Consider linear regression with prior $p(w_d|\lambda_d) = \mathcal{N}(0, \lambda_d^{-1})$ on each weight
- Let's treat precision λ_d as unknown and use a gamma (shape = a , rate = b) prior on it

$$p(\lambda_d) = \text{Gamma}(a, b) = \frac{b^a}{\Gamma(a)} \lambda_d^{a-1} \exp(-b\lambda_d)$$

- Marginalizing the precision leads to a **Student-t prior** on each w_d

$$p(w_d) = \int p(w_d|\lambda_d)p(\lambda_d)d\lambda_d = \frac{b^a \Gamma(a + 1/2)}{\sqrt{2\pi} \Gamma(a)} (b + w_d^2/2)^{-(a+1/2)}$$



- Note: Can make the prior an **uninformative prior** by setting a and b to be very small (e.g., 10^{-4})
- Note: Some other priors on λ_d (e.g., exponential distribution) also result in sparse priors on w_d

Bayesian Linear Regression with Sparse Prior on Weights

- Posterior inference for \mathbf{w} not straightforward since $p(\mathbf{w}) = \prod_{d=1}^D p(w_d)$ is no longer Gaussian
- Approximate inference is usually needed for inferring the full posterior



Bayesian Linear Regression with Sparse Prior on Weights

- Posterior inference for \mathbf{w} not straightforward since $p(\mathbf{w}) = \prod_{d=1}^D p(w_d)$ is no longer Gaussian
- Approximate inference is usually needed for inferring the full posterior
- Many approaches exist (which we will see later)



Bayesian Linear Regression with Sparse Prior on Weights

- Posterior inference for \mathbf{w} not straightforward since $p(\mathbf{w}) = \prod_{d=1}^D p(w_d)$ is no longer Gaussian
- Approximate inference is usually needed for inferring the full posterior
- Many approaches exist (which we will see later)
- Such approaches are mostly in form of alternating estimation of \mathbf{w} and λ



Bayesian Linear Regression with Sparse Prior on Weights

- Posterior inference for \mathbf{w} not straightforward since $p(\mathbf{w}) = \prod_{d=1}^D p(w_d)$ is no longer Gaussian
- Approximate inference is usually needed for inferring the full posterior
- Many approaches exist (which we will see later)
- Such approaches are mostly in form of alternating estimation of \mathbf{w} and λ
 - Estimate λ_d given w_d , estimate w_d given λ_d



Bayesian Linear Regression with Sparse Prior on Weights

- Posterior inference for \mathbf{w} not straightforward since $p(\mathbf{w}) = \prod_{d=1}^D p(w_d)$ is no longer Gaussian
- Approximate inference is usually needed for inferring the full posterior
- Many approaches exist (which we will see later)
- Such approaches are mostly in form of alternating estimation of \mathbf{w} and λ
 - Estimate λ_d given w_d , estimate w_d given λ_d
- Popular approaches: EM, Gibbs sampling, variational inference, etc



Bayesian Linear Regression with Sparse Prior on Weights

- Posterior inference for \mathbf{w} not straightforward since $p(\mathbf{w}) = \prod_{d=1}^D p(w_d)$ is no longer Gaussian
- Approximate inference is usually needed for inferring the full posterior
- Many approaches exist (which we will see later)
- Such approaches are mostly in form of alternating estimation of \mathbf{w} and λ
 - Estimate λ_d given w_d , estimate w_d given λ_d
- Popular approaches: EM, Gibbs sampling, variational inference, etc
- Working with such sparse priors is known as [Sparse Bayesian Learning](#)



Bayesian Linear Regression with Sparse Prior on Weights

- Posterior inference for \mathbf{w} not straightforward since $p(\mathbf{w}) = \prod_{d=1}^D p(w_d)$ is no longer Gaussian
- Approximate inference is usually needed for inferring the full posterior
- Many approaches exist (which we will see later)
- Such approaches are mostly in form of alternating estimation of \mathbf{w} and λ
 - Estimate λ_d given w_d , estimate w_d given λ_d
- Popular approaches: EM, Gibbs sampling, variational inference, etc
- Working with such sparse priors is known as [Sparse Bayesian Learning](#)
 - Used in many models where we want to have sparsity in the weights (very few non-zero weights)



Bayesian Linear Regression with Sparse Prior on Weights

- Posterior inference for \mathbf{w} not straightforward since $p(\mathbf{w}) = \prod_{d=1}^D p(w_d)$ is no longer Gaussian
- Approximate inference is usually needed for inferring the full posterior
- Many approaches exist (which we will see later)
- Such approaches are mostly in form of alternating estimation of \mathbf{w} and λ
 - Estimate λ_d given w_d , estimate w_d given λ_d
- Popular approaches: EM, Gibbs sampling, variational inference, etc
- Working with such sparse priors is known as [Sparse Bayesian Learning](#)
 - Used in many models where we want to have sparsity in the weights (very few non-zero weights)
- Note: We will later look at other ways of getting sparsity (e.g., [spike-and-slab priors](#) defined by binary switch variables for each weight)



Bayesian Logistic Regression

(..a simple, single-parameter, yet **non-conjugate** model)



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn $p(y|\mathbf{x})$: Discriminative Classification and Generative Classification



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn $p(y|\mathbf{x})$: Discriminative Classification and Generative Classification
- **Discriminative Classification:** Model and learn $p(y|\mathbf{x})$ directly



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn $p(y|\mathbf{x})$: Discriminative Classification and Generative Classification
- **Discriminative Classification:** Model and learn $p(y|\mathbf{x})$ directly
 - This approach does not model the distribution of the inputs \mathbf{x}



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn $p(y|\mathbf{x})$: Discriminative Classification and Generative Classification
- **Discriminative Classification:** Model and learn $p(y|\mathbf{x})$ directly
 - This approach does not model the distribution of the inputs \mathbf{x}
- **Generative Classification:** Model and learn $p(y|\mathbf{x})$ “indirectly”



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn $p(y|\mathbf{x})$: Discriminative Classification and Generative Classification
- **Discriminative Classification:** Model and learn $p(y|\mathbf{x})$ directly
 - This approach does not model the distribution of the inputs \mathbf{x}
- **Generative Classification:** Model and learn $p(y|\mathbf{x})$ “indirectly” as $p(y|\mathbf{x}) = \frac{p(y)p(\mathbf{x}|y)}{p(\mathbf{x})}$



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn $p(y|\mathbf{x})$: Discriminative Classification and Generative Classification
- **Discriminative Classification:** Model and learn $p(y|\mathbf{x})$ directly
 - This approach does not model the distribution of the inputs \mathbf{x}
- **Generative Classification:** Model and learn $p(y|\mathbf{x})$ “indirectly” as $p(y|\mathbf{x}) = \frac{p(y)p(\mathbf{x}|y)}{p(\mathbf{x})}$
 - Called generative because, via $p(\mathbf{x}|y)$, we model how the inputs \mathbf{x} of each class are generated



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn $p(y|\mathbf{x})$: Discriminative Classification and Generative Classification
- **Discriminative Classification:** Model and learn $p(y|\mathbf{x})$ directly
 - This approach does not model the distribution of the inputs \mathbf{x}
- **Generative Classification:** Model and learn $p(y|\mathbf{x})$ “indirectly” as $p(y|\mathbf{x}) = \frac{p(y)p(\mathbf{x}|y)}{p(\mathbf{x})}$
 - Called generative because, via $p(\mathbf{x}|y)$, we model how the inputs \mathbf{x} of each class are generated
 - The approach requires first learning **class-marginal** $p(y)$ and **class-conditional** distributions $p(\mathbf{x}|y)$



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn $p(y|\mathbf{x})$: Discriminative Classification and Generative Classification
- **Discriminative Classification:** Model and learn $p(y|\mathbf{x})$ directly
 - This approach does not model the distribution of the inputs \mathbf{x}
- **Generative Classification:** Model and learn $p(y|\mathbf{x})$ “indirectly” as $p(y|\mathbf{x}) = \frac{p(y)p(\mathbf{x}|y)}{p(\mathbf{x})}$
 - Called generative because, via $p(\mathbf{x}|y)$, we model how the inputs \mathbf{x} of each class are generated
 - The approach requires first learning **class-marginal** $p(y)$ and **class-conditional** distributions $p(\mathbf{x}|y)$
 - Usually harder to learn than discriminative but also has some advantages (more on this later)



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn $p(y|\mathbf{x})$: Discriminative Classification and Generative Classification
- **Discriminative Classification:** Model and learn $p(y|\mathbf{x})$ directly
 - This approach does not model the distribution of the inputs \mathbf{x}
- **Generative Classification:** Model and learn $p(y|\mathbf{x})$ “indirectly” as $p(y|\mathbf{x}) = \frac{p(y)p(\mathbf{x}|y)}{p(\mathbf{x})}$
 - Called generative because, via $p(\mathbf{x}|y)$, we model how the inputs \mathbf{x} of each class are generated
 - The approach requires first learning **class-marginal** $p(y)$ and **class-conditional** distributions $p(\mathbf{x}|y)$
 - Usually harder to learn than discriminative but also has some advantages (more on this later)
- Both approaches can be given a non-Bayesian or Bayesian treatment



Probabilistic Models for Classification

- The goal is to learn $p(y|\mathbf{x})$. Here $p(y|\mathbf{x})$ will be a discrete distribution (e.g., Bernoulli, multinoulli)
- Usually two approaches to learn $p(y|\mathbf{x})$: Discriminative Classification and Generative Classification
- **Discriminative Classification:** Model and learn $p(y|\mathbf{x})$ directly
 - This approach does not model the distribution of the inputs \mathbf{x}
- **Generative Classification:** Model and learn $p(y|\mathbf{x})$ “indirectly” as $p(y|\mathbf{x}) = \frac{p(y)p(\mathbf{x}|y)}{p(\mathbf{x})}$
 - Called generative because, via $p(\mathbf{x}|y)$, we model how the inputs \mathbf{x} of each class are generated
 - The approach requires first learning **class-marginal** $p(y)$ and **class-conditional** distributions $p(\mathbf{x}|y)$
 - Usually harder to learn than discriminative but also has some advantages (more on this later)
- Both approaches can be given a non-Bayesian or Bayesian treatment
 - The Bayesian treatment won't rely on point estimates but infer the posterior over unknowns



Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$



Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**



Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w})$$



Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu$$



Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x})$$



Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})}$$



Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$



Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

where $\mathbf{w} \in \mathbb{R}^D$ is the weight vector.



Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

where $\mathbf{w} \in \mathbb{R}^D$ is the weight vector. Also note that $p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \mu$

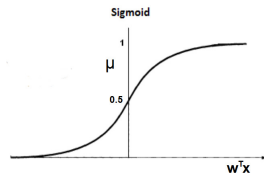


Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

where $\mathbf{w} \in \mathbb{R}^D$ is the weight vector. Also note that $p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \mu$

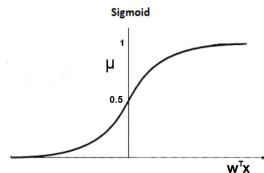


Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

where $\mathbf{w} \in \mathbb{R}^D$ is the weight vector. Also note that $p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \mu$



- A large positive (negative) “score” $\mathbf{w}^\top \mathbf{x}$ means large probability of label being 1 (0)

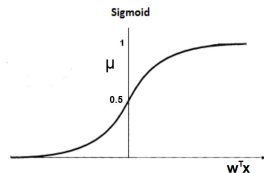


Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

where $\mathbf{w} \in \mathbb{R}^D$ is the weight vector. Also note that $p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \mu$



- A large positive (negative) “score” $\mathbf{w}^\top \mathbf{x}$ means large probability of label being 1 (0)
- Is sigmoid the only way to convert the score into a probability?

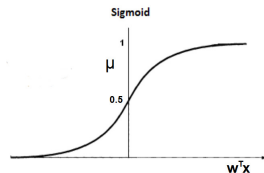


Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

where $\mathbf{w} \in \mathbb{R}^D$ is the weight vector. Also note that $p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \mu$



- A large positive (negative) “score” $\mathbf{w}^\top \mathbf{x}$ means large probability of label being 1 (0)
- Is sigmoid the only way to convert the score into a probability?
 - No, while LR does that, there exist models that define μ in other ways

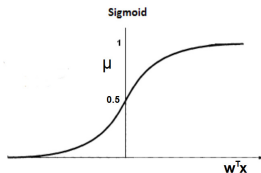


Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

where $\mathbf{w} \in \mathbb{R}^D$ is the weight vector. Also note that $p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \mu$



- A large positive (negative) “score” $\mathbf{w}^\top \mathbf{x}$ means large probability of label being 1 (0)
- Is sigmoid the only way to convert the score into a probability?
 - No, while LR does that, there exist models that define μ in other ways. E.g. **Probit Regression**

$$\mu = p(y = 1|\mathbf{x}, \mathbf{w}) = \Phi(\mathbf{w}^\top \mathbf{x})$$

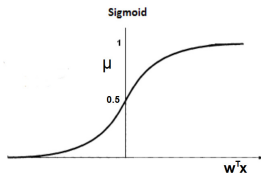


Discriminative Classification via Logistic Regression

- **Logistic Regression** (LR) is an example of discriminative **binary** classification, i.e., $y \in \{0, 1\}$
- Logistic Regression models \mathbf{x} to y relationship using the **sigmoid function**

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

where $\mathbf{w} \in \mathbb{R}^D$ is the weight vector. Also note that $p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \mu$



- A large positive (negative) “score” $\mathbf{w}^\top \mathbf{x}$ means large probability of label being 1 (0)
- Is sigmoid the only way to convert the score into a probability?
 - No, while LR does that, there exist models that define μ in other ways. E.g. **Probit Regression**

$$\mu = p(y = 1|\mathbf{x}, \mathbf{w}) = \Phi(\mathbf{w}^\top \mathbf{x}) \quad (\text{where } \Phi \text{ denotes the CDF of } \mathcal{N}(0, 1))$$



Logistic Regression

- The LR classification rule is

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$



Logistic Regression

- The LR classification rule is

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$



Logistic Regression

- The LR classification rule is

$$p(y = 1|x, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0|x, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- This implies a **Bernoulli likelihood** model for the labels

$$p(y|x, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x}))$$



Logistic Regression

- The LR classification rule is

$$p(y = 1|x, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0|x, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- This implies a **Bernoulli likelihood** model for the labels

$$p(y|x, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$



Logistic Regression

- The LR classification rule is

$$p(y = 1|x, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0|x, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- This implies a **Bernoulli likelihood** model for the labels

$$p(y|x, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, we can do point estimation for \mathbf{w} by maximizing the log-likelihood (or minimizing the **negative log-likelihood**).



Logistic Regression

- The LR classification rule is

$$p(y = 1|x, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0|x, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- This implies a **Bernoulli likelihood** model for the labels

$$p(y|x, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, we can do point estimation for \mathbf{w} by maximizing the log-likelihood (or minimizing the **negative log-likelihood**). This is basically MLE.

$$\mathbf{w}_{MLE} = \arg \max_{\mathbf{w}} \sum_{n=1}^N \log p(y_n|x_n, \mathbf{w})$$



Logistic Regression

- The LR classification rule is

$$p(y = 1|x, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0|x, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- This implies a **Bernoulli likelihood** model for the labels

$$p(y|x, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, we can do point estimation for \mathbf{w} by maximizing the log-likelihood (or minimizing the **negative log-likelihood**). This is basically MLE.

$$\mathbf{w}_{MLE} = \arg \max_{\mathbf{w}} \sum_{n=1}^N \log p(y_n | \mathbf{x}_n, \mathbf{w}) = \arg \min_{\mathbf{w}} - \sum_{n=1}^N \log p(y_n | \mathbf{x}_n, \mathbf{w})$$



Logistic Regression

- The LR classification rule is

$$p(y = 1|x, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0|x, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- This implies a **Bernoulli likelihood** model for the labels

$$p(y|x, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, we can do point estimation for \mathbf{w} by maximizing the log-likelihood (or minimizing the **negative log-likelihood**). This is basically MLE.

$$\mathbf{w}_{MLE} = \arg \max_{\mathbf{w}} \sum_{n=1}^N \log p(y_n|x_n, \mathbf{w}) = \arg \min_{\mathbf{w}} - \sum_{n=1}^N \log p(y_n|x_n, \mathbf{w}) = \arg \min_{\mathbf{w}} \text{NLL}(\mathbf{w})$$



Logistic Regression

- The LR classification rule is

$$p(y = 1|x, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0|x, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- This implies a **Bernoulli likelihood** model for the labels

$$p(y|x, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, we can do point estimation for \mathbf{w} by maximizing the log-likelihood (or minimizing the **negative log-likelihood**). This is basically MLE.

$$\mathbf{w}_{MLE} = \arg \max_{\mathbf{w}} \sum_{n=1}^N \log p(y_n|x_n, \mathbf{w}) = \arg \min_{\mathbf{w}} - \sum_{n=1}^N \log p(y_n|x_n, \mathbf{w}) = \arg \min_{\mathbf{w}} \text{NLL}(\mathbf{w})$$

- Convex loss function. Global minima. Both first order and second order methods widely used.



Logistic Regression

- The LR classification rule is

$$p(y = 1|x, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0|x, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- This implies a **Bernoulli likelihood** model for the labels

$$p(y|x, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, we can do point estimation for \mathbf{w} by maximizing the log-likelihood (or minimizing the **negative log-likelihood**). This is basically MLE.

$$\mathbf{w}_{MLE} = \arg \max_{\mathbf{w}} \sum_{n=1}^N \log p(y_n|x_n, \mathbf{w}) = \arg \min_{\mathbf{w}} - \sum_{n=1}^N \log p(y_n|x_n, \mathbf{w}) = \arg \min_{\mathbf{w}} \text{NLL}(\mathbf{w})$$

- Convex loss function. Global minima. Both first order and second order methods widely used.

- Can also add a regularizer on \mathbf{w} to prevent overfitting. This corresponds to doing MAP estimation with a prior on \mathbf{w} , i.e., $\mathbf{w}_{MAP} = \arg \max_{\mathbf{w}} [\sum_{n=1}^N \log p(y_n|x_n, \mathbf{w}) + \log p(\mathbf{w})]$



Bayesian Logistic Regression

- MLE/MAP only gives a point estimate. We would like to infer the full posterior over \mathbf{w}



Bayesian Logistic Regression

- MLE/MAP only gives a point estimate. We would like to infer the full posterior over \mathbf{w}
- Recall that the **likelihood model** is Bernoulli

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$



Bayesian Logistic Regression

- MLE/MAP only gives a point estimate. We would like to infer the full posterior over \mathbf{w}
- Recall that the **likelihood model** is Bernoulli

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Just like the Bayesian linear regression case, let's use a Gaussian **prior** on \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I}_D) \propto \exp\left(-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}\right)$$



Bayesian Logistic Regression

- MLE/MAP only gives a point estimate. We would like to infer the full posterior over \mathbf{w}
- Recall that the **likelihood model** is Bernoulli

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Just like the Bayesian linear regression case, let's use a Gaussian **prior** on \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I}_D) \propto \exp\left(-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}\right)$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, where \mathbf{X} is $N \times D$ and \mathbf{y} is $N \times 1$, the posterior over \mathbf{w}

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$$



Bayesian Logistic Regression

- MLE/MAP only gives a point estimate. We would like to infer the full posterior over \mathbf{w}
- Recall that the **likelihood model** is Bernoulli

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Just like the Bayesian linear regression case, let's use a Gaussian **prior** on \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I}_D) \propto \exp\left(-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}\right)$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, where \mathbf{X} is $N \times D$ and \mathbf{y} is $N \times 1$, the posterior over \mathbf{w}

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} = \frac{\prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})}{\int \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$$



Bayesian Logistic Regression

- MLE/MAP only gives a point estimate. We would like to infer the full posterior over \mathbf{w}
- Recall that the **likelihood model** is Bernoulli

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Just like the Bayesian linear regression case, let's use a Gaussian **prior** on \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I}_D) \propto \exp\left(-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}\right)$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, where \mathbf{X} is $N \times D$ and \mathbf{y} is $N \times 1$, the posterior over \mathbf{w}

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} = \frac{\prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})}{\int \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$$

- The denominator is intractable in general (logistic-Bernoulli and Gaussian are not conjugate)



Bayesian Logistic Regression

- MLE/MAP only gives a point estimate. We would like to infer the full posterior over \mathbf{w}
- Recall that the **likelihood model** is Bernoulli

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Just like the Bayesian linear regression case, let's use a Gaussian **prior** on \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I}_D) \propto \exp\left(-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}\right)$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, where \mathbf{X} is $N \times D$ and \mathbf{y} is $N \times 1$, the posterior over \mathbf{w}

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} = \frac{\prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})}{\int \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$$

- The denominator is intractable in general (logistic-Bernoulli and Gaussian are not conjugate)
 - Can't get a closed form expression for $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$. Must approximate it!



Bayesian Logistic Regression

- MLE/MAP only gives a point estimate. We would like to infer the full posterior over \mathbf{w}
- Recall that the **likelihood model** is Bernoulli

$$p(y|\mathbf{x}, \mathbf{w}) = \text{Bernoulli}(\sigma(\mathbf{w}^\top \mathbf{x})) = \left[\frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^y \left[\frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})} \right]^{(1-y)}$$

- Just like the Bayesian linear regression case, let's use a Gaussian **prior** on \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(0, \lambda^{-1} \mathbf{I}_D) \propto \exp\left(-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}\right)$$

- Given N observations $(\mathbf{X}, \mathbf{y}) = \{\mathbf{x}_n, y_n\}_{n=1}^N$, where \mathbf{X} is $N \times D$ and \mathbf{y} is $N \times 1$, the posterior over \mathbf{w}

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} = \frac{\prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})}{\int \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$$

- The denominator is intractable in general (logistic-Bernoulli and Gaussian are not conjugate)
 - Can't get a closed form expression for $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$. Must approximate it!
 - Several ways to do it, e.g., MCMC, variational inference, **Laplace approximation** (next class)



Next Class

- Laplace approximation
- Computing posterior and posterior predictive for logistic regression
- Properties/benefits of Bayesian logistic regression
- Bayesian approach to generative classification

