

Probabilistic Linear Regression

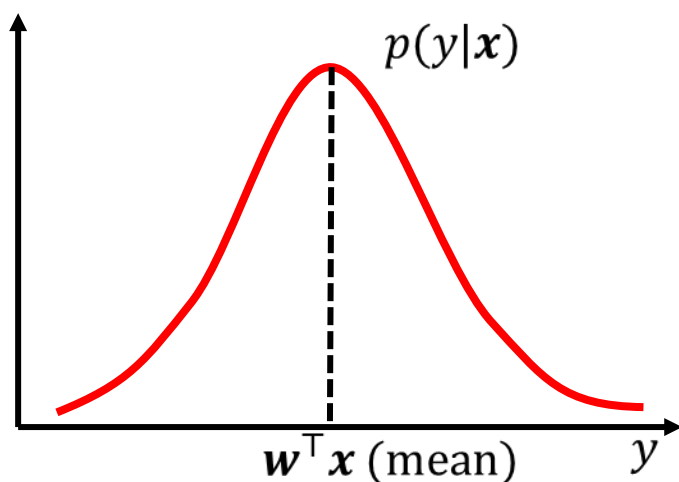
CS771: Introduction to Machine Learning

Nisheeth

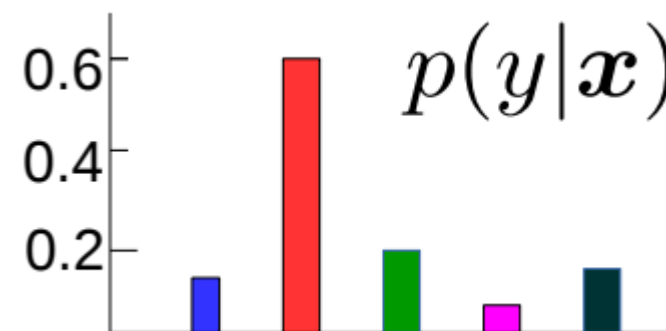
Probabilistic Models for Supervised Learning

- Goal: Learn the conditional distribution of output given input, i.e., $p(y|\mathbf{x})$

Probabilistic Linear Regression



Probabilistic Classification



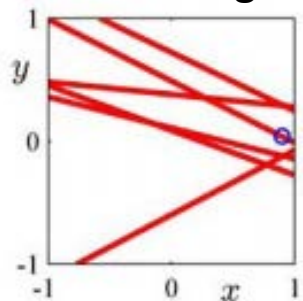
- $p(y|\mathbf{x})$ is more informative than a single prediction y
 - From $p(y|\mathbf{x})$, can get “expected” or “most likely” output y
 - For classifn, “soft” predictions (e.g., rather than yes/no, prob. of “yes”)
 - “Uncertainty” in the predicted output y (e.g., by looking at the variance of $p(y|\mathbf{x})$)
- Can also learn a distribution over the model params using **fully Bayesian inference**



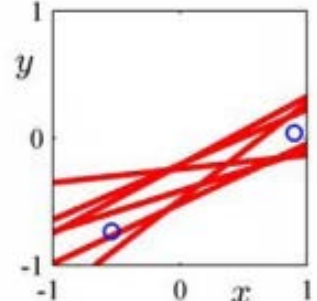
Distribution over model parameters

- Recall that linear/ridge regression gave a single “optimal” weight vector
- With a probabilistic model for linear regression, we have two options
 - Use MLE/MAP to get a single “optimal” weight vector
 - Use fully Bayesian inference to learn a distribution over weight vectors (figure below)

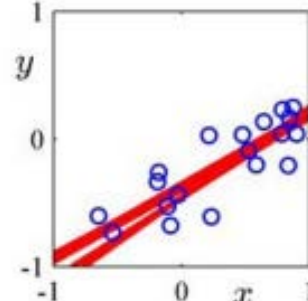
One training ex



Two training ex



A few more training ex



Rather than returning just a single “best” solution (a line in this example), the fully Bayesian approach would give us several “probable” lines (consistent with training data) by learning the **full posterior distribution over the model parameters** (each of which corresponds to a line)



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

Posterior predictive distribution by doing posterior weighted averaging over all possible \mathbf{w} , not just the most likely one. Thus more robust predictions especially if we are uncertain about the best solution.

Predictive distribution using a single \mathbf{w} (plug-in predictive distribution)

How important/like this \mathbf{w} is under the posterior distribution (its posterior probability)



Probabilistic Models for Supervised Learning

- Usually two ways to model the conditional distribution $p(y|\mathbf{x})$
- Approach 1: Don't model \mathbf{x} , and model $p(y|\mathbf{x})$ directly using a prob. distribution

“discriminative” sup learning

Gaussian distribution

Probabilistic linear regression

$$p(y|\mathbf{x}, \mathbf{w}) = \mathcal{N}(y|\mathbf{w}^\top \mathbf{x}, \beta^{-1})$$

The “sigmoid” function

Probabilistic linear binary classification

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

We assume the conditional distribution to be some appropriate distribution and treat the weights \mathbf{w} as learnable parameters of the model (using MLE/MAP/fully Bayesian inference). Need not be a linear model – can replace $\mathbf{w}^\top \mathbf{x}$ by a nonlinear function $f(\mathbf{x})$



- Approach 2: Model both \mathbf{x} and y via their joint distr. and get the conditional as

“generative” sup learning

Here θ denotes all the model parameters that we need to model the joint distribution of \mathbf{x} and y (will see examples later)

$$p(y|\mathbf{x}, \theta) = \frac{p(\mathbf{x}, y|\theta)}{p(\mathbf{x}|\theta)}$$

Prob. distribution of inputs from class k

$$p(y = k|\mathbf{x}, \theta) = \frac{p(\mathbf{x}, y=k|\theta)}{p(\mathbf{x}|\theta)} = \frac{p(\mathbf{x}|y = k, \theta)p(y=k|\theta)}{\sum_{\ell=1}^K p(\mathbf{x}|y = \ell, \theta)p(y=\ell|\theta)}$$

For a multi-class classification model with K classes

Called “generative” because we are learning the generative distributions for output as well as inputs



Brief Detour (Gaussian Distribution)

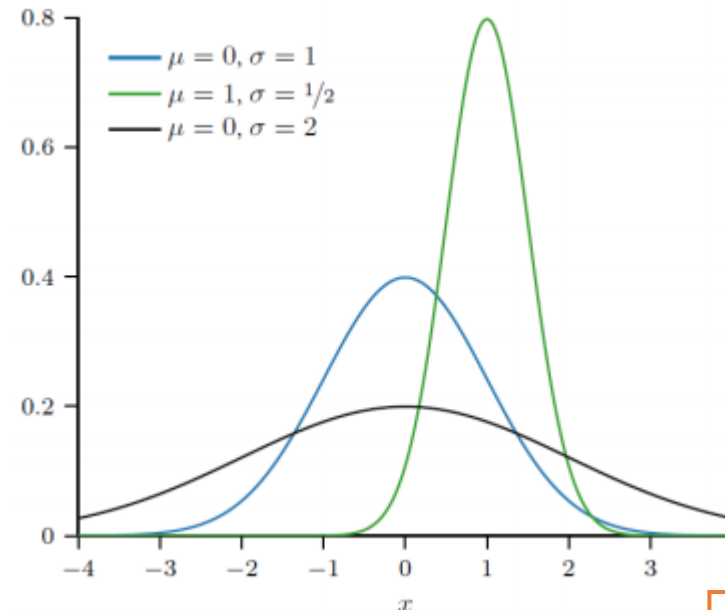


Gaussian Distribution (Univariate)

- Distribution over real-valued scalar random variables $x \in \mathbb{R}$
- Defined by a scalar mean μ and a scalar variance σ^2

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right]$$

- Mean: $\mathbb{E}[x] = \mu$
- Variance: $\text{var}[x] = \sigma^2$
- Inverse of variance is called **precision**: $\beta = \frac{1}{\sigma^2}$.



Gaussian PDF in terms of precision

$$\mathcal{N}(x|\mu, \beta) = \sqrt{\frac{\beta}{2\pi}} \exp\left[-\frac{\beta}{2}(x - \mu)^2\right]$$

Gaussian Distribution (Multivariate)

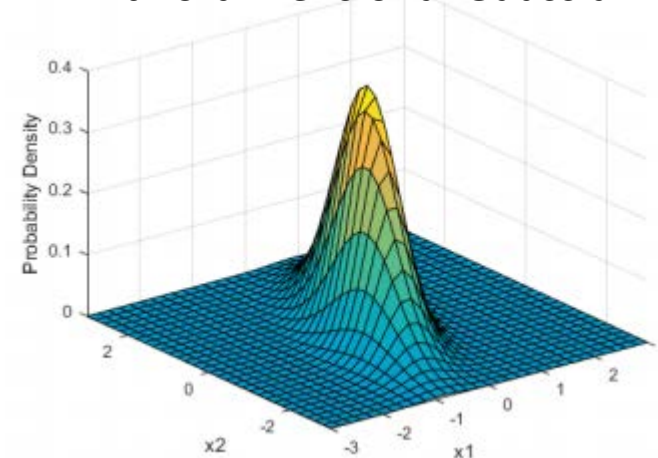
- Distribution over real-valued vector random variables $\mathbf{x} \in \mathbb{R}^D$
- Defined by a mean vector $\boldsymbol{\mu} \in \mathbb{R}^D$ and a covariance matrix $\boldsymbol{\Sigma}$

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp[-(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})]$$

- Note: The cov. matrix $\boldsymbol{\Sigma}$ must be symmetric and PSD
 - All eigenvalues are positive
 - $\mathbf{z}^\top \boldsymbol{\Sigma} \mathbf{z} \geq 0$ for any real vector \mathbf{z}

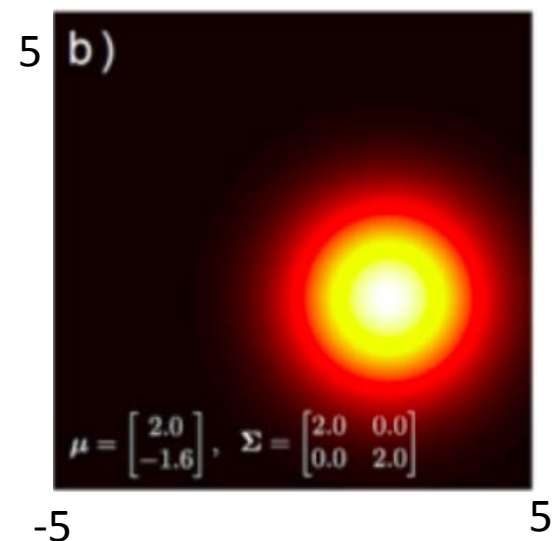
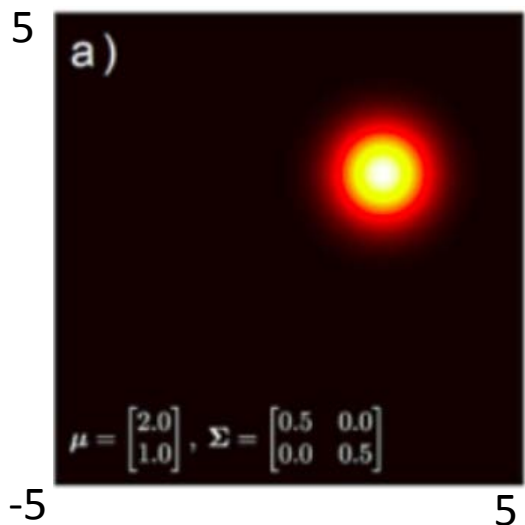
- The covariance matrix also controls the shape of the Gaussian

A two-dimensional Gaussian

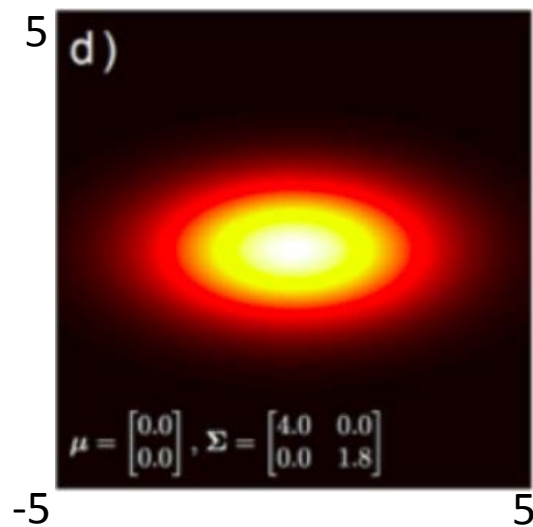
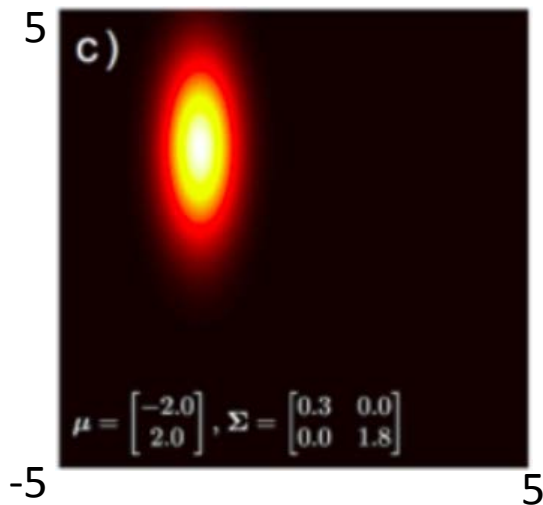


Covariance Matrix for Multivariate Gaussian

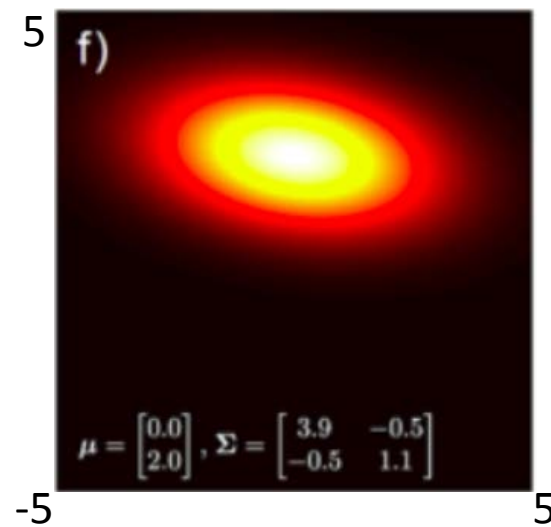
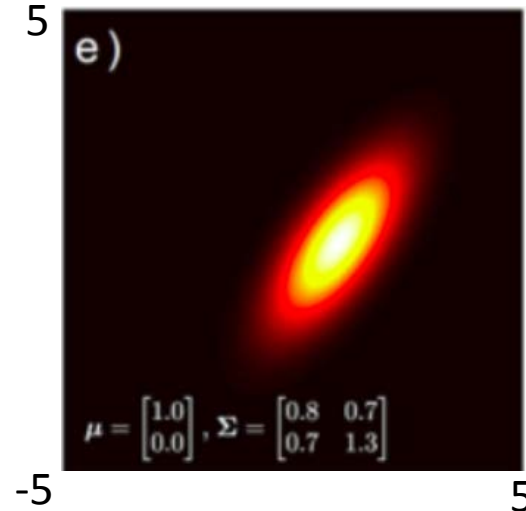
Spherical Covariance



Diagonal Covariance



Full Covariance



Probabilistic Linear Regression

$$p(y|\mathbf{x}, \mathbf{w}) = \mathcal{N}(y|\mathbf{w}^\top \mathbf{x}, \beta^{-1})$$

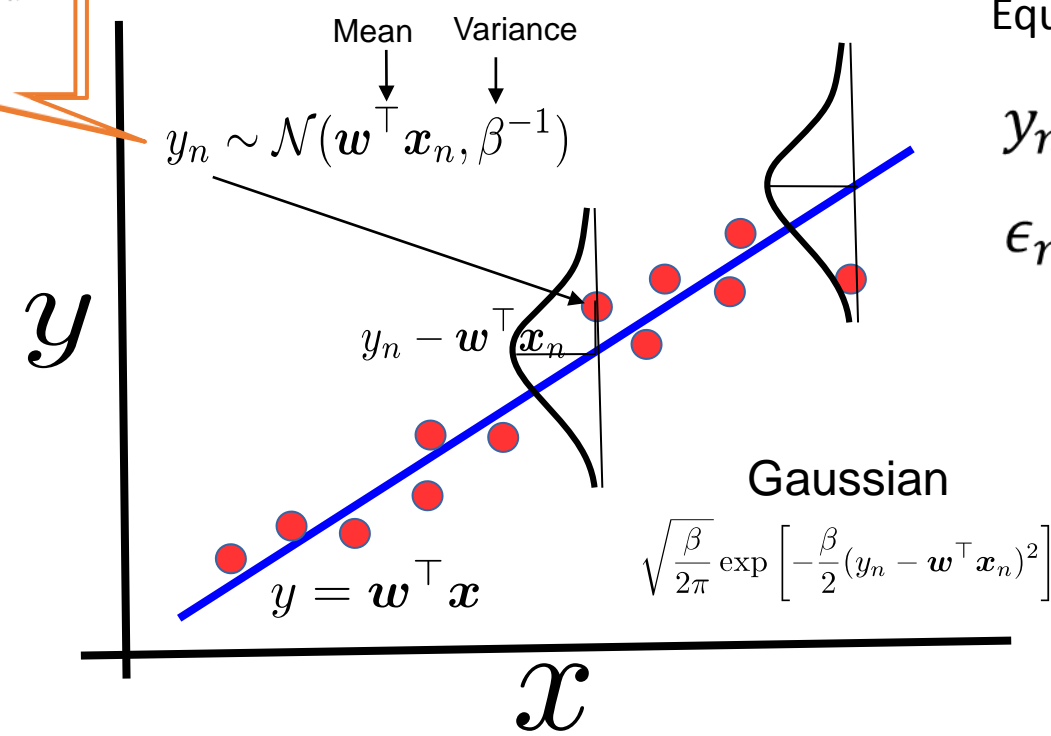
[Nice tutorial](#)



Linear Regression: A Probabilistic View

Defines our likelihood model:
 $p(y_n | \mathbf{w}, \mathbf{x}_n)$ - Gaussian

Output y_n assumed generated from a Gaussian with mean $\mathbf{w}^\top \mathbf{x}_n$



Output y_n generated from a linear model and then zero mean Gaussian noise added

Equivalently:

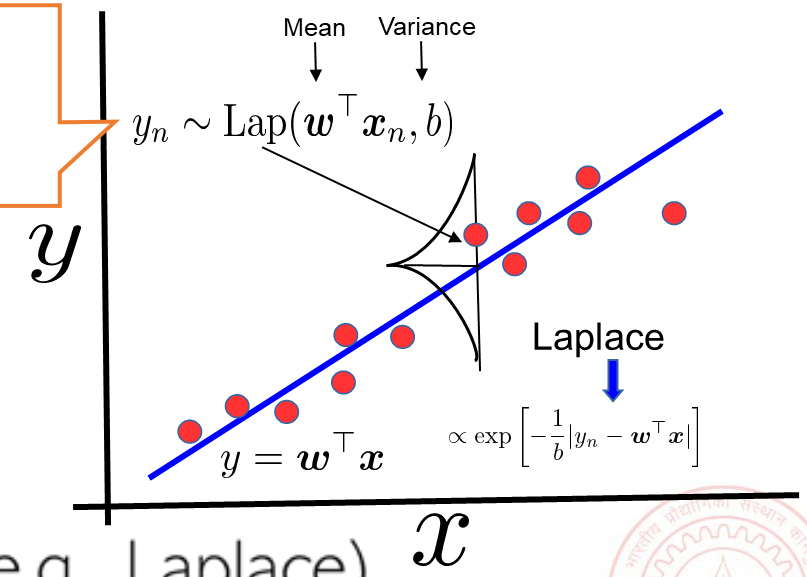
$$y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$$

$$\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$$

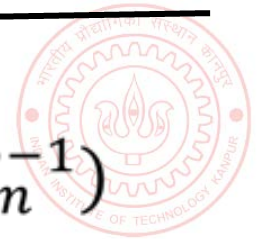
Note the term in the Gaussian's exponent – just like a squared error we saw for least squares regression



Using a Laplace distribution would correspond to using an absolute loss



- Several variants of this basic model are possible
 - Other distributions to model the additive noise (e.g., Laplace)
 - Different noise variance/precision for each output: $y_n \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta_n^{-1})$



MLE for Probabilistic Linear Regression

- Since each likelihood term is a Gaussian, we have

Also note that \mathbf{x}_n is fixed here but the likelihood depend on it, so it is being conditioned on

$$p(y_n | \mathbf{w}, \mathbf{x}_n) = \mathcal{N}(y_n | \mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp \left[-\frac{\beta}{2} (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 \right]$$

Exercise: Verify that you can also write the overall likelihood as a single N dimensional Gaussian with mean $\mathbf{X}\mathbf{w}$ and cov. matrix $\beta^{-1}\mathbf{I}_N$

- Thus the overall likelihood (assuming i.i.d. responses) will be

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n | \mathbf{x}_n, \mathbf{w}) = \left(\frac{\beta}{2\pi} \right)^{N/2} \exp \left[-\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 \right]$$

- Log-likelihood (ignoring constants w.r.t. \mathbf{w})

$$\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) \propto -\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

MLE for probabilistic linear regression with Gaussian noise is equivalent to least squares regression without any regularization (with solution $\hat{\mathbf{w}}_{MLE} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$)



- Negative log likelihood (NLL) in this case is similar to **squared loss function**

MAP Estimation for Prob. Lin. Reg.: The Prior

- For MAP estimation, we need a prior distribution over the parameters $\mathbf{w} \in \mathbb{R}^D$
- A reasonable prior for real-valued vectors can be a multivariate Gaussian

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_0, \Sigma)$$

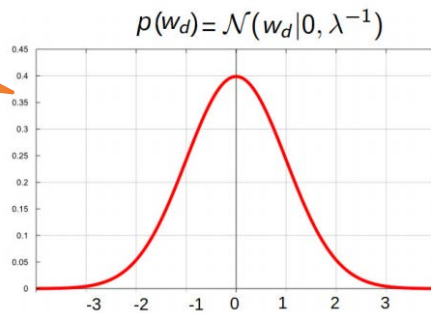
Equivalent to saying that *a priori* we expect the solution to be close to some vector \mathbf{w}_0
(subject to Σ being such that the variances is not too large)

- A specific example of a multivariate Gaussian prior in this problem

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda^{-1} \mathbf{I}_D) = \prod_{d=1}^D \mathcal{N}(w_d | 0, \lambda^{-1}) = \prod_{d=1}^D p(w_d)$$

Omitting λ for brevity

The precision λ of the Gaussian prior controls how aggressively the prior pushes the elements towards mean (0)



This is essentially like a regularizer that pushes elements of \mathbf{w} to be small (we will see shortly)

Equivalent to saying that *a priori* we expect each element of the solution to be close to 0 (i.e., “small”)

$$\mathcal{N}(w_d | 0, \lambda^{-1}) = \sqrt{\frac{\lambda}{2\pi}} \exp\left[-\frac{\lambda}{2} w_d^2\right]$$

$$\mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda^{-1} \mathbf{I}_D) = \left(\frac{\lambda}{2\pi}\right)^{D/2} \exp\left[-\frac{\lambda}{2} \sum_{d=1}^D w_d^2\right] = \left(\frac{\lambda}{2\pi}\right)^{D/2} \exp\left[-\frac{\lambda}{2} \mathbf{w}^T \mathbf{w}\right]$$



MAP Estimation for Probabilistic Linear Regression¹³

- The MAP objective (log-posterior) will be the **log-likelihood** + **$\log p(\mathbf{w})$**

$$-\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 - \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

In the likelihood and prior, ignored terms that don't depend on \mathbf{w}

- Maximizing this is equivalent to minimizing the following w.r.t. \mathbf{w}

$$\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{\beta} \mathbf{w}^\top \mathbf{w}$$

Not surprising since MAP estimation indeed optimizes a regularized loss function!



- This is equivalent to ridge regression with regularization hyperparameter $\frac{\lambda}{\beta}$

- The solution will be $\hat{\mathbf{w}}_{MAP} = (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$



Fully Bayesian Inference for Prob. Linear Regression¹⁴

- Can also compute the full posterior distribution over \mathbf{w}

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

For brevity, we have not shown the dependence of the various distributions here on the hyperparameters λ and β

- Likelihood and prior are conjugate (both Gaussians) - posterior will be Gaussian

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$
$$\boldsymbol{\mu}_N = (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$$
$$\boldsymbol{\Sigma}_N = (\beta \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1}$$

Posterior's mean is the same as the MAP solution since the mean and mode of a Gaussian are the same

Note: λ and β are assumed to be fixed; otherwise, the problem is a bit harder (beyond the scope of CS771)

We now have a distribution over the possible solutions – it has a mean but we can generate other plausible solutions by sampling from this posterior. Each sample will give a weight vector



Prob. Linear Regression: The Predictive Distribution ¹⁵

- Want the predictive distribution $p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y})$ of the output y_* for a new input \mathbf{x}_*
- With MLE/MAP estimate of \mathbf{w} , we will use the plug-in predictive

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) = \mathcal{N}(\mathbf{w}_{MLE}^\top \mathbf{x}_*, \beta^{-1}) \quad \text{- MLE prediction}$$
$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx p(y_* | \mathbf{x}_*, \mathbf{w}_{MAP}) = \mathcal{N}(\mathbf{w}_{MAP}^\top \mathbf{x}_*, \beta^{-1}) \quad \text{- MAP prediction}$$

- When doing fully Bayesian inference, can compute the **posterior predictive dist.**

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

Not true in general for Prob. Lin. Reg. but because the hyperparameters λ and β are treated as fixed

- Requires an integral but has a closed form

Mean prediction

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\mu_N^\top \mathbf{x}_*, \beta^{-1} + \mathbf{x}_*^\top \Sigma_N \mathbf{x}_*)$$

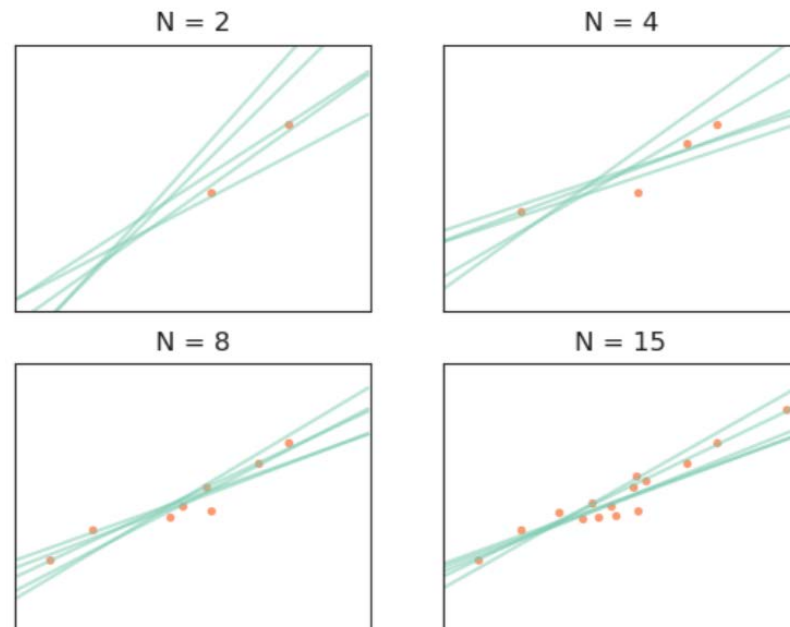
Input-specific predictive variance unlike the MLE/MAP based predictive where it was β^{-1} (and was same for all test inputs)

- Input-specific predictive uncertainty useful in problems where we want confidence estimates of the predictions made by the model (e.g., Active Learning)

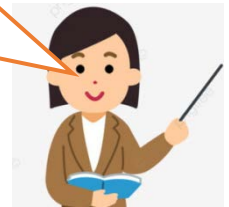


Fully Bayesian Linear Regression – Pictorially

- Each sample from posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$ will give a weight vector \mathbf{w}
 - In case of lin. reg., each weight vector corresponds to a regression line



The posterior sort of represents an ensemble of solutions (not all are equally good but we can use all of them in an “importance-weighted” fashion to make the prediction using the posterior predictive distribution)



Importance of each solution in this ensemble is its posterior probability $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$

- Each weight vector will give a different set of predictions on test data
 - These different predictions will give us a variance (uncertainty) estimate in model’s prediction
 - The uncertainty decreases as N increases (we become more sure when we see more training data)



MLE, MAP/Fully Bayesian Lin. Reg: Summary

- MLE/MAP give point estimate of \mathbf{w}
 - MLE/MAP based prediction uses that single point estimate of \mathbf{w}
- Fully Bayesian approach gives the full posterior of \mathbf{w}
 - Fully Bayesian prediction does posterior averaging (computes posterior predictive distribution)
- Some things to keep in mind:
 - MLE estimation of a parameter leads to **unregularized solutions**
 - MAP estimation of a parameter leads to **regularized solutions**
 - A **Gaussian likelihood** model corresponds to using **squared loss**
 - A **Gaussian prior** on parameters acts as an ℓ_2 regularizer
 - Other likelihoods/priors can be chosen (result in other loss functions and regularizers)

E.g., using Laplace distribution for likelihood is equivalent to absolute loss, using it as a prior is equivalent to ℓ_1 regularization



Evaluation Measures for Regression Models

- Plotting the prediction \hat{y}_n vs truth y_n for the validation/test set
- Residual Sum of Squares (RSS) on the validation/test set

$$RSS(\mathbf{w}) = \sum_{n=1}^N (y_n - \hat{y}_n)^2$$

- RMSE (Root Mean Squared Error) $\triangleq \sqrt{\frac{1}{N} RSS(\mathbf{w})}$
- Coefficient of determination or R^2

$$R^2 = 1 - \frac{\sum_{n=1}^N (y_n - \hat{y}_n)^2}{\sum_{n=1}^N (y_n - \bar{y})^2}$$

Unlike RSS and RMSE, it is always between 0 and 1 and hence interpretable

\bar{y} is empirical mean of true responses, i.e., $\frac{1}{N} \sum_{n=1}^N y_n$

"relative" error w.r.t. a model that makes a constant prediction \bar{y} for all inputs

Plots of true vs predicted outputs and R^2 for two regression models

