

# Parameter Estimation and Prediction in Probabilistic Models

CS772A: Probabilistic Machine Learning

Piyush Rai

# Parameter Estimation: Summary of approaches

- Usually one of the following approaches taken
  - A single “best” **point estimate** of the parameters by optimizing an objective function

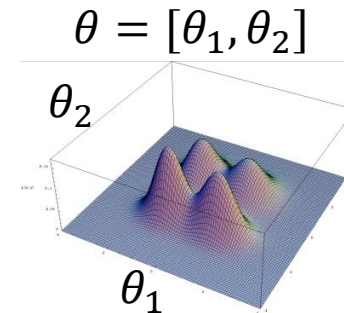
$$\hat{\theta} = \operatorname{argmax}_{\theta} f(\mathcal{D}; \theta)$$

$f$  can be log-likelihood (for MLE) or log-posterior (for MAP)

- A distribution over the parameters (conditioned on observed data  $\mathcal{D}$ )

The posterior distribution

$$p(\theta|\mathcal{D})$$

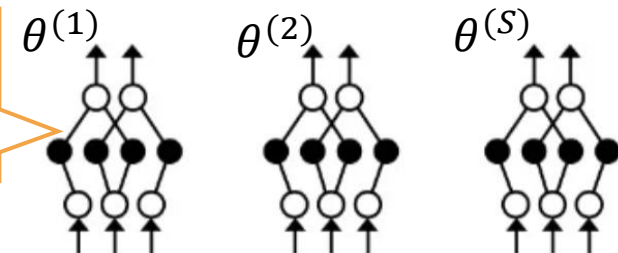


- A set/ensemble of point estimates of parameters (applying approach 1 multiple times)

Computing multiple point estimates, each using a different subset of the training data

$$\{\hat{\theta}(\mathcal{D}') : \mathcal{D}' \sim p^*\}$$

Ensemble (e.g., training a deep neural net with multiple different initializations)



# Making Predictions: The Predictive Distribution

- Prediction: Computing  $p(\mathcal{D}_*|\mathcal{D})$  where  $\mathcal{D}_*$  is test data and  $\mathcal{D}$  is training data
- If we only have a single point estimate  $\hat{\theta}$  then

$$p(\mathcal{D}_*|\mathcal{D}) \approx p(\mathcal{D}_*|\hat{\theta})$$

- If we have computed  $p(\theta|\mathcal{D})$  then the predictive distribution can be defined as

$$\begin{aligned} p(\mathcal{D}_*|\mathcal{D}) &= \int p(\mathcal{D}_*, \theta|\mathcal{D}) d\theta = \int p(\mathcal{D}_*|\theta, \mathcal{D})p(\theta|\mathcal{D}) d\theta \\ &= \int p(\mathcal{D}_*|\theta)p(\theta|\mathcal{D}) d\theta \\ &= \mathbb{E}_{p(\theta|\mathcal{D})}[p(\mathcal{D}_*|\theta)] \end{aligned}$$

If the proper expectation w.r.t. the posterior is computed, as in this case,  $p(\mathcal{D}_*|\mathcal{D})$  is called the **posterior predictive distribution (PPD)**

- If we don't have  $p(\theta|\mathcal{D})$  but a set/ensemble of estimates  $\{\theta^{(i)}\}_{i=1}^S$  then

$$p(\mathcal{D}_*|\mathcal{D}) \approx \frac{1}{S} \sum_{i=1}^S p(\mathcal{D}_*|\theta^{(i)})$$



# Marginal Likelihood

- Recall the posterior distribution

$$p(\theta|\mathcal{D}, \alpha) = \frac{p(\mathcal{D}, \theta|\alpha)}{p(\mathcal{D}|\alpha)} = \frac{p(\mathcal{D}|\theta)p(\theta|\alpha)}{p(\mathcal{D}|\alpha)}$$

- Marginal likelihood  $p(\mathcal{D}|\alpha)$  is an important quantity and is by definition

$$p(\mathcal{D}|\alpha) = \int p(\mathcal{D}|\theta)p(\theta|\alpha)d\theta = \mathbb{E}_{p(\theta|\alpha)}[p(\mathcal{D}|\theta)]$$

- Can think of it as the likelihood averaged over all  $\theta$ 's from the prior
- Useful quantity in general. For example, we can find the best hyperparameter  $\alpha$  as

$$\hat{\alpha} = \operatorname{argmax}_{\alpha} \log p(\mathcal{D}|\alpha) \quad \text{No validation set required! ☺}$$

- Hard to compute however (because of the expectation involved)
  - But approximations can be made



# A “Shortcut”: PPD using Marginal Likelihood

- PPD, by definition, is obtained by the following marginalization

$$p(\mathcal{D}_*|\mathcal{D}) = \int p(\mathcal{D}_*|\theta)p(\theta|\mathcal{D}) d\theta$$

- Can also compute PPD without computing the posterior! Some ways:

1. Using a ratio of marginal likelihoods as follows

Follows simply from Bayes rule

$$p(a|b) = \frac{p(a,b)}{p(b)}$$

$$p(\mathcal{D}_*|\mathcal{D}) = \frac{p(\mathcal{D}_*, \mathcal{D})}{p(\mathcal{D})}$$

Joint marginal likelihood  
for training and test data

Marginal likelihood for  
training data

2. If  $p(\mathcal{D}_*|\mathcal{D})$  can be obtained easily from the joint  $p(\mathcal{D}_*, \mathcal{D})$

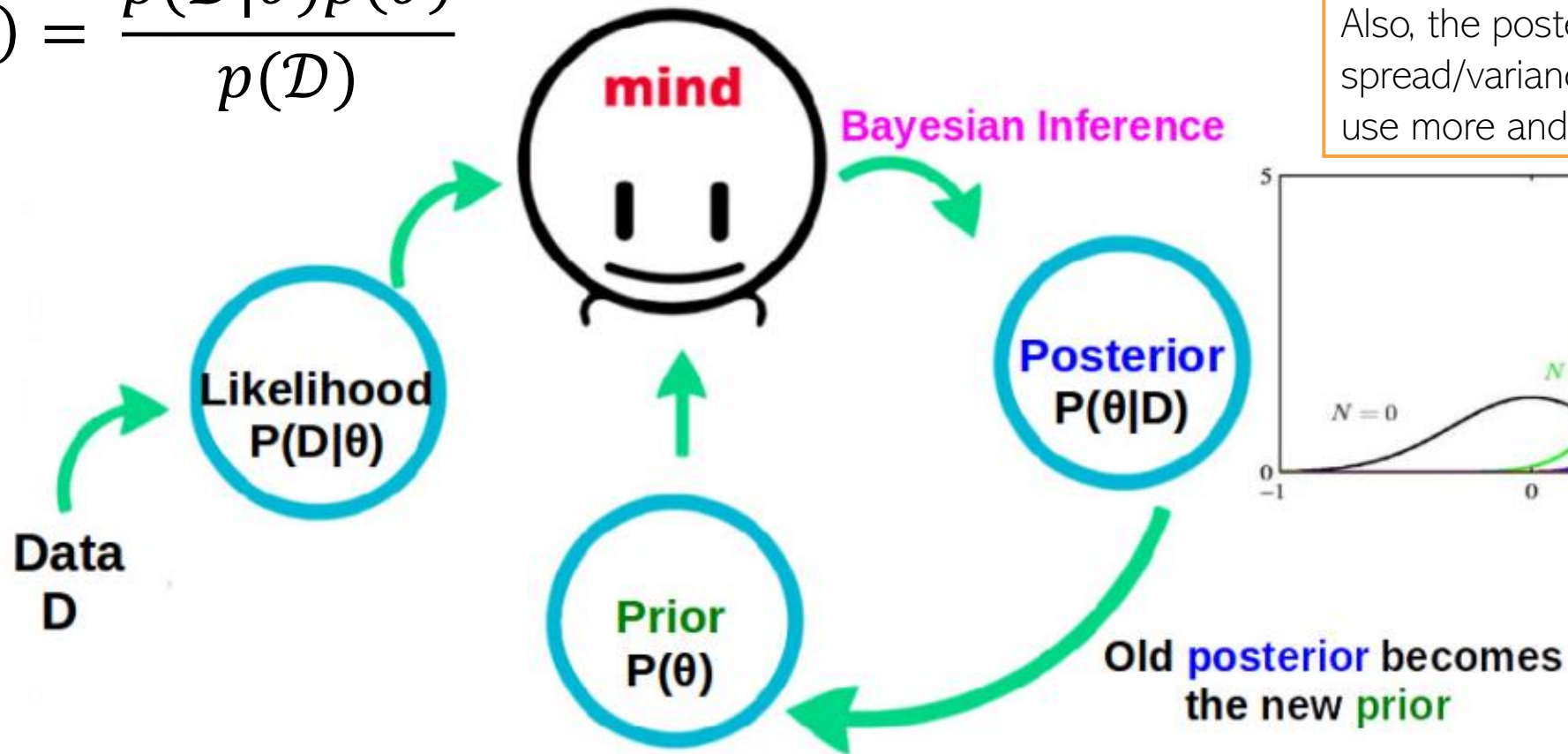
- Note that the PPD  $p(\mathcal{D}_*|\mathcal{D})$  is also a conditional distribution
- For some distributions (e.g., Gaussian), conditionals can be easily derived from joint

Will see this being used we we  
study Gaussian Process (GP)

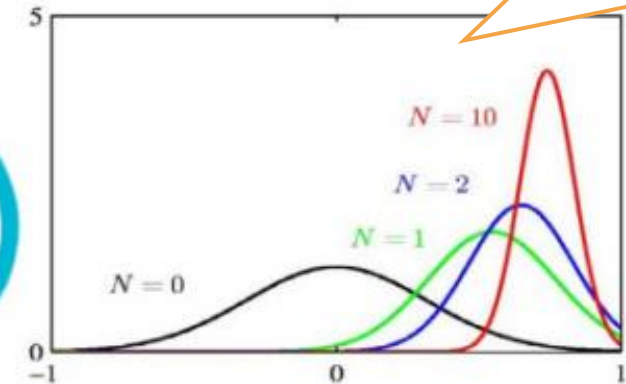
# An Important Aspect: Posterior Updates

- Posterior updates in Bayesian inference can naturally be done in an online fashion

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$



Also, the posterior's spread/variance gets smaller as we use more and more data to infer it



# Coming up..

- Parameter estimation (point est. and posterior) and predictive distribution for
  - Bernoulli observation model (binary-valued observations)
  - Multinoulli observation model (discrete-valued observations)
  - Gaussian observation model (real-valued observations)



# The IID Assumption

- Assume that, conditioned on  $\theta$ , observations are independently and identically distributed (i.i.d. assumption). Depending on the problem, this may look like:

Supervised generative model  
(both inputs and output are modeled using a distribution)

$$(\mathbf{x}_n, y_n) \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}, y | \theta) \quad \longrightarrow \quad p(\mathcal{D} | \theta) = \prod_{i=1}^N p(\mathbf{x}_i, y_i | \theta)$$

Supervised discriminative model  
(only the output is modeled using a distribution); input is assumed "given" and not modeled

$$y_n \stackrel{\text{i.i.d.}}{\sim} p(y | \mathbf{x}, \theta) \quad \longrightarrow \quad p(\mathcal{D} | \theta) = \prod_{i=1}^N p(y_i | \mathbf{x}_i, \theta)$$

Unsupervised generative model (there are only inputs; no labels)

$$\mathbf{x}_n \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x} | \theta) \quad \longrightarrow \quad p(\mathcal{D} | \theta) = \prod_{i=1}^N p(\mathbf{x}_i | \theta)$$

- Assume that both training and test data come from the same distribution
  - This assumption, although standard, may be violated in real-world applications of ML and there are "adaptation" methods to handle that





# Bernoulli Observation Model



# Estimating a Coin's Bias

- Consider a sequence of  $N$  coin toss outcomes (observations)
- Each observation  $y_n$  is a binary **random variable**. Head:  $y_n = 1$ , Tail:  $y_n = 0$
- Each  $y_n$  is assumed generated by a **Bernoulli distribution** with param  $\theta \in (0,1)$

Probability  
of a head

Likelihood or  
observation model

$$p(y_n|\theta) = \text{Bernoulli}(y_n|\theta) = \theta^{y_n} (1 - \theta)^{1-y_n}$$

- Here  $\theta$  the unknown param (probability of head). Let's do MLE
- Log-likelihood:  $\sum_{n=1}^N \log p(y_n|\theta) = \sum_{n=1}^N [y_n \log \theta + (1 - y_n) \log (1 - \theta)]$
- Maximizing log-lik, or minimizing neg. log-lik (NLL) w.r.t.  $\theta$  gives

assuming i.i.d. data

I tossed a coin 5 times – gave 1 head and 4 tails. Does it mean  $\theta = 0.2$ ?? The MLE approach says so. What if I see 0 head and 5 tails. Does it mean  $\theta = 0$ ?

$$\theta_{MLE} = \frac{\sum_{n=1}^N y_n}{N}$$

Thus MLE solution is simply the fraction of heads! 😊 Makes intuitive sense!

Indeed, with a small number of training observations, MLE may overfit and may not be reliable. An alternative is MAP estimation which can incorporate a **prior distribution** over  $\theta$

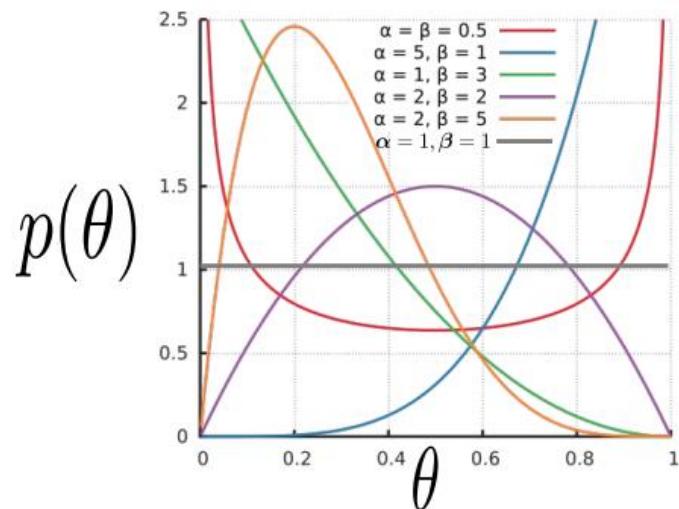


# Estimating a Coin's Bias

- Let's do MAP estimation for the bias of the coin
- Each likelihood term is Bernoulli

$$p(y_n|\theta) = \text{Bernoulli}(y_n|\theta) = \theta^{y_n} (1 - \theta)^{1-y_n}$$

- Also need a prior since we want to do MAP estimation
- Since  $\theta \in (0,1)$ , a reasonable choice of prior for  $\theta$  would be [Beta distribution](#)



$$p(\theta|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

The gamma function

Using  $\alpha = 1$  and  $\beta = 1$  will make the Beta prior a uniform prior

$\alpha$  and  $\beta$  (both non-negative reals) are the two hyperparameters of this Beta prior

Can set these based on intuition, cross-validation, or even learn them

# Estimating a Coin's Bias

- The **log posterior** for the coin-toss model is log-lik + log-prior

$$LP(\theta) = \sum_{n=1}^N \log p(y_n|\theta) + \log p(\theta|\alpha, \beta)$$

- Plugging in the expressions for Bernoulli and Beta and ignoring any terms that don't depend on  $\theta$ , the log posterior simplifies to

$$LP(\theta) = \sum_{n=1}^N [y_n \log \theta + (1 - y_n) \log(1 - \theta)] + (\alpha - 1) \log \theta + (\beta - 1) \log(1 - \theta)$$

- Maximizing the above log post. (or min. of its negative) w.r.t.  $\theta$  gives

Using  $\alpha = 1$  and  $\beta = 1$  gives us the same solution as MLE

Recall that  $\alpha = 1$  and  $\beta = 1$  for Beta distribution is in fact equivalent to a uniform prior (hence making MAP equivalent to MLE)

$$\theta_{MAP} = \frac{\sum_{n=1}^N y_n + \alpha - 1}{N + \alpha + \beta - 2}$$

Such interpretations of prior's hyperparameters as being "pseudo-observations" exist for various other prior distributions as well (in particular, distributions belonging to "exponential family" of distributions)

Prior's hyperparameters have an interesting interpretation. Can think of  $\alpha - 1$  and  $\beta - 1$  as the number of heads and tails, respectively, before starting the coin-toss experiment (akin to "pseudo-observations")



# The Posterior Distribution

- Let's do fully Bayesian inference and compute the posterior distribution

- Bernoulli likelihood:  $p(y_n|\theta) = \text{Bernoulli}(y_n|\theta) = \theta^{y_n} (1 - \theta)^{1-y_n}$

- Beta prior:  $p(\theta) = \text{Beta}(\theta|\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$

- The posterior can be computed as

Hyperparams  $\alpha, \beta$  not shown for brevity

$$p(\theta|\mathbf{y}) = \frac{p(\theta)p(\mathbf{y}|\theta)}{p(\mathbf{y})} = \frac{p(\theta) \prod_{n=1}^N p(y_n|\theta)}{p(\mathbf{y})} = \frac{\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \prod_{n=1}^N \theta^{y_n} (1-\theta)^{1-y_n}}{\int \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \prod_{n=1}^N \theta^{y_n} (1-\theta)^{1-y_n} d\theta}$$

Number of heads ( $N_1$ )

Number of tails ( $N_0$ )

$\theta^{\sum_{n=1}^N y_n} (1 - \theta)^{N - \sum_{n=1}^N y_n}$

- Here, even without computing the denominator (marg lik), we can identify the posterior

- It is Beta distribution since  $p(\theta|\mathbf{y}) \propto \theta^{\alpha+N_1-1} (1 - \theta)^{\beta+N_0-1}$

- Thus  $p(\theta|\mathbf{y}) = \text{Beta}(\theta|\alpha + N_1, \beta + N_0)$

Hint: Use the fact that the posterior must integrate to 1  
 $\int p(\theta|\mathbf{y}) d\theta = 1$

Exercise: Show that the normalization constant equals

$$\frac{\Gamma(\alpha + \beta + N)}{\Gamma(\alpha + \sum_{n=1}^N y_n) \Gamma(\beta + N - \sum_{n=1}^N y_n)}$$



- Here, finding the posterior boiled down to simply “multiply, add stuff, and identify”

- Here, posterior has the same form as prior (both Beta): property of **conjugate priors**.

# Conjugacy and Conjugate Priors

- Many pairs of distributions are conjugate to each other
  - Bernoulli (likelihood) + Beta (prior)  $\Rightarrow$  Beta posterior
  - Binomial (likelihood) + Beta (prior)  $\Rightarrow$  Beta posterior
  - Multinomial (likelihood) + Dirichlet (prior)  $\Rightarrow$  Dirichlet posterior
  - Poisson (likelihood) + Gamma (prior)  $\Rightarrow$  Gamma posterior
  - Gaussian (likelihood) + Gaussian (prior)  $\Rightarrow$  Gaussian posterior
  - and many other such pairs ..

Not true in general, but in some cases (e.g., the variance of the Gaussian likelihood is fixed)

- Tip: If two distr are conjugate to each other, their functional forms are similar

- Example: Bernoulli and Beta have the forms

$$\text{Bernoulli}(y|\theta) = \theta^y (1 - \theta)^{1-y}$$

$$\text{Beta}(\theta|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

This is why, when we multiply them while computing the posterior, the exponents get added and we get the same form for the posterior as the prior but with just updated hyperparameter. Also, we can identify the posterior and its hyperparameters simply by inspection

- More on conjugate priors when we look at **exponential family** distributions



# Predictive Distribution

- Suppose we want to compute the prob that the next outcome  $\mathbf{y}_{N+1}$  will be head (=1)
- The **posterior predictive distribution** (averaging over all  $\theta$ 's weighted by their respective posterior probabilities)

$$\begin{aligned}
 p(y_{N+1} = 1|\mathbf{y}) &= \int_0^1 p(y_{N+1} = 1, \theta|\mathbf{y}) d\theta = \int_0^1 p(y_{N+1} = 1|\theta)p(\theta|\mathbf{y}) d\theta \\
 &= \int_0^1 \theta \times p(\theta|\mathbf{y}) d\theta \\
 &= \mathbb{E}_{p(\theta|\mathbf{y})}[\theta] \\
 &= \frac{\alpha + N_1}{\alpha + \beta + N}
 \end{aligned}$$

Expectation of  $\theta$  w.r.t. the Beta posterior distribution  $p(\theta|\mathbf{y}) = \text{Beta}(\theta|\alpha + N_1, \beta + N_0)$

For models where likelihood and prior are conjugate to each other, the PPD can be computed easily in closed form (more on this when we talk about **exponential family** distributions)

- Therefore the PPD will be

$$p(y_{N+1}|\mathbf{y}) = \text{Bernoulli}(y_{N+1}|\mathbb{E}_{p(\theta|\mathbf{y})}[\theta])$$

- The **plug-in predictive** distribution using a point estimate  $\hat{\theta}$  (e.g., using MLE/MAP)

$$p(y_{N+1} = 1|\mathbf{y}) \approx p(y_{N+1} = 1|\hat{\theta}) = \hat{\theta} \longrightarrow p(y_{N+1}|\mathbf{y}) = \text{Bernoulli}(y_{N+1}|\hat{\theta})$$



# Multinoulli Observation Model





# The Posterior Distribution

MLE/MAP left as an exercise

17

- Assume  $N$  discrete obs  $\mathbf{y} = \{y_1, y_2, \dots, y_N\}$  with each  $y_n \in \{1, 2, \dots, K\}$ , e.g.,
  - $y_n$  represents the outcome of a dice roll with  $K$  faces
  - $y_n$  represents the class label of the  $n^{th}$  example in a classification problem (total  $K$  classes)
  - $y_n$  represents the identity of the  $n^{th}$  word in a sequence of words

- Assume **likelihood** to be multinoulli with unknown params  $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_K]$

$$p(y_n|\boldsymbol{\pi}) = \text{multinoulli}(y_n|\boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{\mathbb{I}[y_n=k]}$$

Generalization of Bernoulli to  $K > 2$  discrete outcomes

- $\boldsymbol{\pi}$  is a vector of probabilities (“probability vector”), e.g.,
  - Biases of the  $K$  sides of the dice
  - Prior class probabilities in multi-class classification ( $p(y_n = k) = \pi_k$ )
  - Probabilities of observing each word of the  $K$  words in a vocabulary

Called the **concentration parameter** of the Dirichlet (assumed known for now)

Large values of  $\boldsymbol{\alpha}$  will give a Dirichlet peaked around its mean (next slides illustrates this)

Each  $\alpha_k \geq 0$

- Assume a **conjugate prior** (Dirichlet) on  $\boldsymbol{\pi}$  with hyperparams  $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \dots, \alpha_K]$

$$p(\boldsymbol{\pi}|\boldsymbol{\alpha}) = \text{Dirichlet}(\boldsymbol{\pi}|\alpha_1, \dots, \alpha_K) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k-1} = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \pi_k^{\alpha_k-1}$$

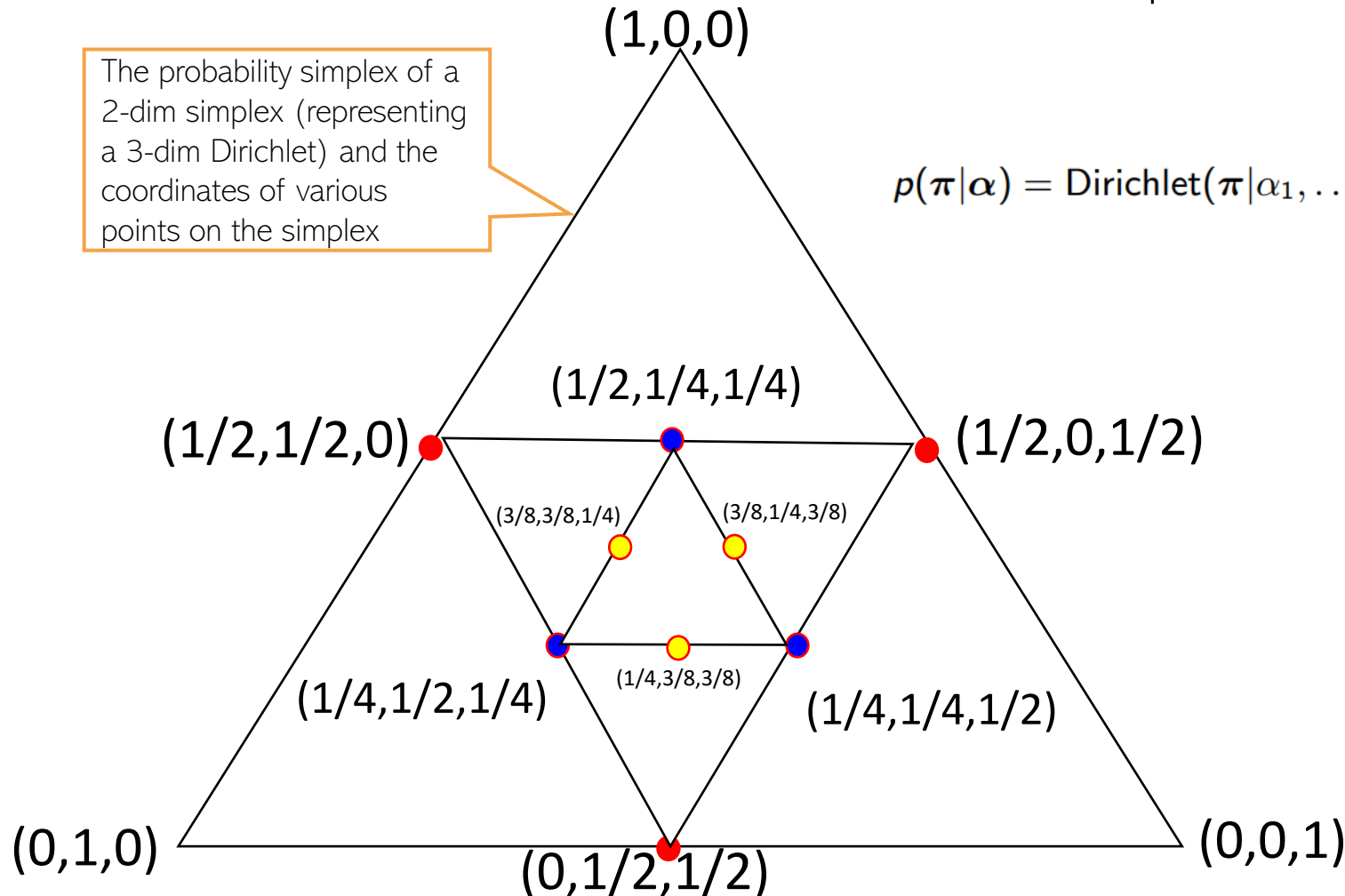
Generalization of Beta to  $K$ -dimensional **probability vectors**

# Brief Detour: Dirichlet Distribution

Basically, probability vectors

- An important distribution. Models non-neg. vectors  $\boldsymbol{\pi}$  that also sum to one
- A random draw from  $K$ -dim Dirich. will be a point under  $(K-1)$ -dim probability simplex

The probability simplex of a 2-dim simplex (representing a 3-dim Dirichlet) and the coordinates of various points on the simplex

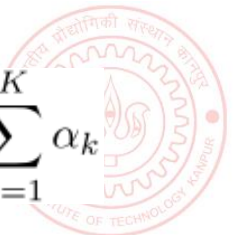


$$p(\boldsymbol{\pi}|\boldsymbol{\alpha}) = \text{Dirichlet}(\boldsymbol{\pi}|\alpha_1, \dots, \alpha_K) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k-1} = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k=1}^K \pi_k^{\alpha_k-1}$$

$$\text{Mean} = \left[ \frac{\alpha_1}{\sum_{k=1}^K \alpha_k}, \dots, \frac{\alpha_K}{\sum_{k=1}^K \alpha_k} \right]$$

$$\text{Mode} = \left[ \frac{\alpha_1 - 1}{\sum_{k=1}^K \alpha_k - K}, \dots, \frac{\alpha_K - 1}{\sum_{k=1}^K \alpha_k - K} \right] (\alpha_k > 1)$$

$$\text{var}(\pi_k) = \frac{\alpha_k(\alpha_0 - \alpha_k)}{\alpha_0^2(\alpha_0 + 1)} \quad \alpha_0 = \sum_{k=1}^K \alpha_k$$



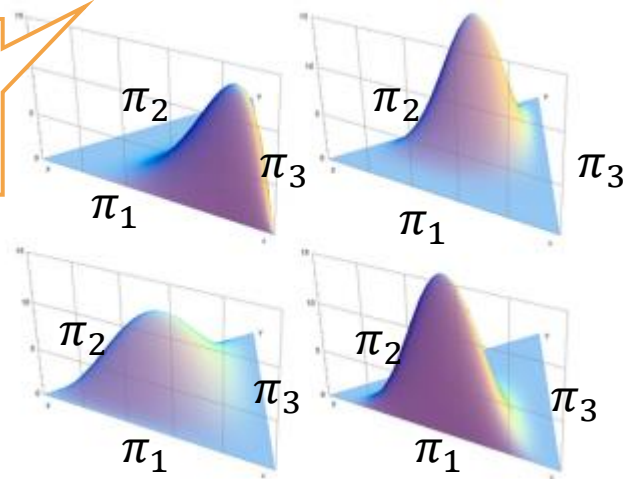
# Brief Detour: Dirichlet Distribution

- A visualization of Dirichlet distribution for different values of concentration param

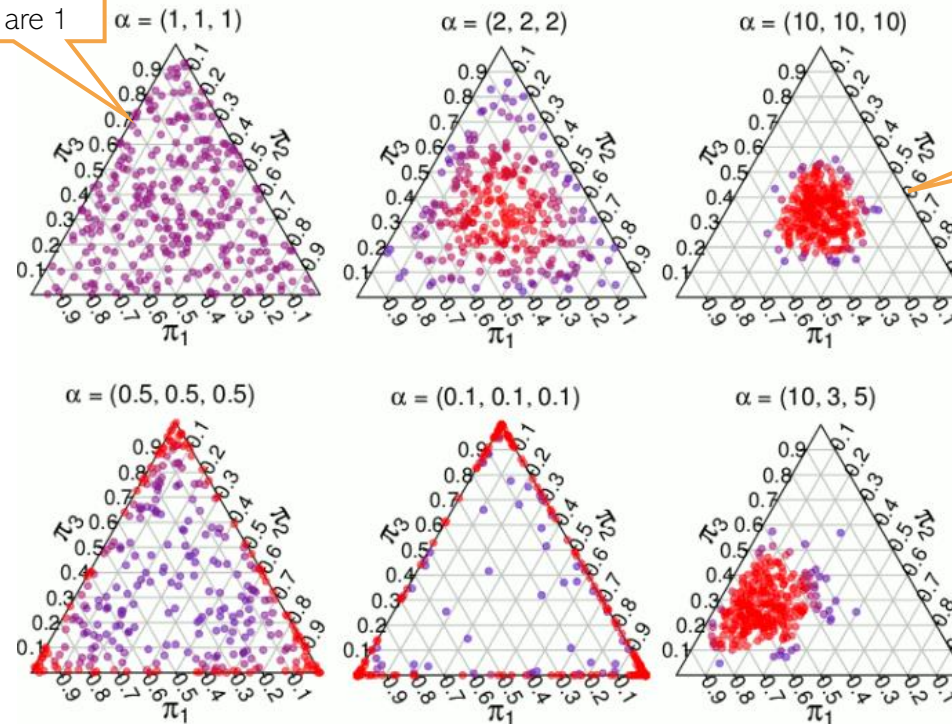
Visualizations of PDFs of some 3-dim Dirichlet distributions (each generated using a different conc. Param vector  $\alpha$ )

Like a uniform distribution if all  $\alpha_k$ 's are 1

$\alpha$  controls the shape of the Dirichlet (just like Beta distribution's hyperparameters)



Draws from a 3-dimensional Dirichlet with different  $\alpha$



All  $\alpha_k$ 's large results in peak around the center of the simplex

- Interesting fact: Can generate a  $K$ -dim Dirichlet random variable by independently generating  $K$  gamma random variables and normalizing them to sum to 1



# The Posterior Distribution

- Posterior  $p(\boldsymbol{\pi}|\mathbf{y})$  is easy to compute due to conjugacy b/w **multinoulli** and **Dir.**

$$p(\boldsymbol{\pi}|\mathbf{y}, \boldsymbol{\alpha}) = \frac{p(\boldsymbol{\pi}, \mathbf{y}|\boldsymbol{\alpha})}{p(\mathbf{y}|\boldsymbol{\alpha})} = \frac{p(\boldsymbol{\pi}|\boldsymbol{\alpha})p(\mathbf{y}|\boldsymbol{\pi}, \boldsymbol{\alpha})}{p(\mathbf{y}|\boldsymbol{\alpha})} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Marg-lik}} = \frac{p(\boldsymbol{\pi}|\boldsymbol{\alpha})p(\mathbf{y}|\boldsymbol{\pi})}{p(\mathbf{y}|\boldsymbol{\alpha})}$$

Don't need to compute for this case because of conjugacy

Marg-lik =  $\int p(\boldsymbol{\pi}|\boldsymbol{\alpha})p(\mathbf{y}|\boldsymbol{\pi})d\boldsymbol{\pi}$

- Assuming  $y_n$ 's are i.i.d. given  $\boldsymbol{\pi}$ ,  $p(\mathbf{y}|\boldsymbol{\pi}) = \prod_{n=1}^N p(y_n|\boldsymbol{\pi})$ , and therefore

$$p(\boldsymbol{\pi}|\mathbf{y}, \boldsymbol{\alpha}) \propto \prod_{k=1}^K \pi_k^{\alpha_k - 1} \times \prod_{n=1}^N \prod_{k=1}^K \pi_k^{\mathbb{I}[y_n=k]} = \prod_{k=1}^K \pi_k^{\alpha_k + \sum_{n=1}^N \mathbb{I}[y_n=k] - 1}$$

- Even without computing marg-lik,  $p(\mathbf{y}|\boldsymbol{\alpha})$ , we can see that the posterior is Dirichlet
- Denoting  $N_k = \sum_{n=1}^N \mathbb{I}[y_n = k]$ , number of observations with value  $k$

$$p(\boldsymbol{\pi}|\mathbf{y}, \boldsymbol{\alpha}) = \text{Dirichlet}(\boldsymbol{\pi}|\alpha_1 + N_1, \alpha_2 + N_2, \dots, \alpha_K + N_K)$$

- Note:  $N_1, N_2, \dots, N_K$  are the sufficient statistics for this estimation problem
  - We only need the suff-stats to estimate the parameters and values of individual observations aren't needed (another property from exponential family of distributions – more on this later)

Similar to number of heads and tails for the coin bias estimation problem



# The Predictive Distribution

- Finally, let's also look at the **posterior predictive distribution** for this model
- PPD is the prob distr of a new  $y_* \in \{1, 2, \dots, K\}$ , given training data  $\mathbf{y} = \{y_1, y_2, \dots, y_N\}$ 

Will be a multinoulli. Just need to estimate the probabilities of each of the  $K$  outcomes

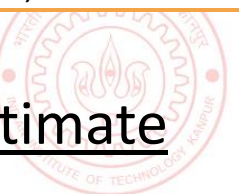
$$p(y_* | \mathbf{y}, \boldsymbol{\alpha}) = \int p(y_* | \boldsymbol{\pi}) p(\boldsymbol{\pi} | \mathbf{y}, \boldsymbol{\alpha}) d\boldsymbol{\pi}$$
- $p(y_* | \boldsymbol{\pi}) = \text{multinoulli}(y_* | \boldsymbol{\pi})$ ,  $p(\boldsymbol{\pi} | \mathbf{y}, \boldsymbol{\alpha}) = \text{Dirichlet}(\boldsymbol{\pi} | \alpha_1 + N_1, \alpha_2 + N_2, \dots, \alpha_K + N_K)$
- Can compute the posterior predictive probability for each of the  $K$  possible outcomes

$$\begin{aligned} p(y_* = k | \mathbf{y}, \boldsymbol{\alpha}) &= \int p(y_* = k | \boldsymbol{\pi}) p(\boldsymbol{\pi} | \mathbf{y}, \boldsymbol{\alpha}) d\boldsymbol{\pi} \\ &= \int \pi_k \times \text{Dirichlet}(\boldsymbol{\pi} | \alpha_1 + N_1, \alpha_2 + N_2, \dots, \alpha_K + N_K) d\boldsymbol{\pi} \\ &= \frac{\alpha_k + N_k}{\sum_{k=1}^K \alpha_k + N} \quad (\text{Expectation of } \pi_k \text{ w.r.t the Dirichlet posterior}) \end{aligned}$$

- Thus PPD is multinoulli with probability vector  $\left\{ \frac{\alpha_k + N_k}{\sum_{k=1}^K \alpha_k + N} \right\}_k^K$ 

Note how these probabilities have been "smoothened" due to the use of the prior + the averaging over the posterior

A similar effect was achieved in the Beta-Bernoulli model, too
- Plug-in predictive will also be multinoulli but with prob vector given by the point estimate of  $\boldsymbol{\pi}$



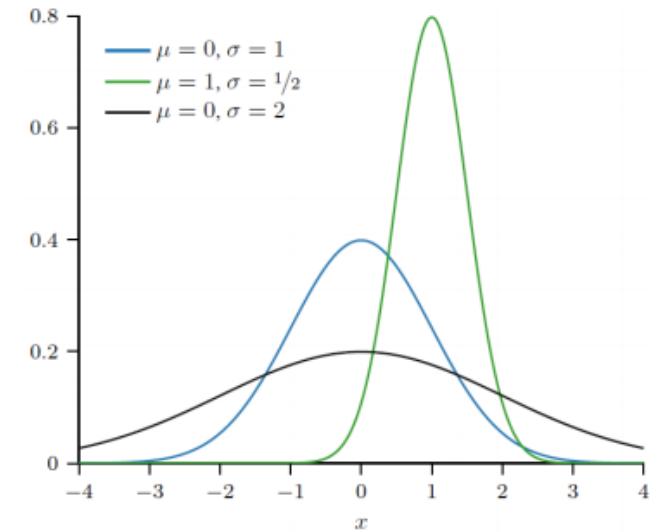
# Gaussian Observation Model



# Gaussian Distribution (Univariate)

- Distribution over real-valued scalar random variables  $X \in \mathbb{R}$ , e.g., height of students in a class
- Defined by a scalar mean  $\mu$  and a scalar variance  $\sigma^2$

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



- 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

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



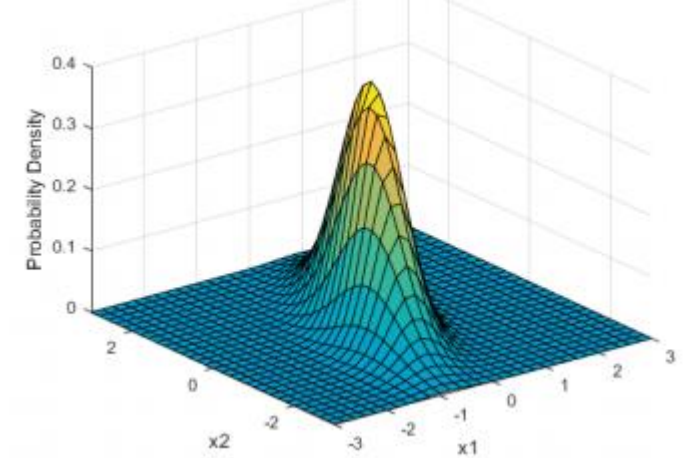
# Gaussian Distribution (Multivariate)

- 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  $\boldsymbol{\Sigma}$

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

- 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





# Linear Gaussian Model (LGM)

- LGM defines a noisy **lin. transform** of a Gaussian r.v.  $\boldsymbol{\theta}$  with  $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$

Both  $\boldsymbol{\theta}$  and  $\mathbf{y}$  are vectors (can be of different sizes)

Also assume  $\mathbf{A}, \mathbf{b}, \boldsymbol{\Lambda}, \mathbf{L}$  to be known; only  $\boldsymbol{\theta}$  is unknown

$$\mathbf{y} = \mathbf{A}\boldsymbol{\theta} + \mathbf{b} + \boldsymbol{\epsilon}$$

Noise vector - independently and drawn from  $\mathcal{N}(\boldsymbol{\epsilon}|\mathbf{0}, \mathbf{L}^{-1})$

- Easy to see that, conditioned on  $\boldsymbol{\theta}$ ,  $\mathbf{y}$  too has a Gaussian distribution

Conditional distribution

$$p(\mathbf{y}|\boldsymbol{\theta}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\theta} + \mathbf{b}, \mathbf{L}^{-1})$$

- Assume  $p(\boldsymbol{\theta})$  as prior and  $p(\mathbf{y}|\boldsymbol{\theta})$  as the likelihood, and defining  $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^\top \mathbf{L} \mathbf{A})^{-1}$

Posterior of  $\boldsymbol{\theta}$

$$p(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})} = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\Sigma}(\mathbf{A}^\top \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}), \boldsymbol{\Sigma})$$

Marginal distribution

$$p(\mathbf{y}) = \int p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^\top + \mathbf{L}^{-1})$$

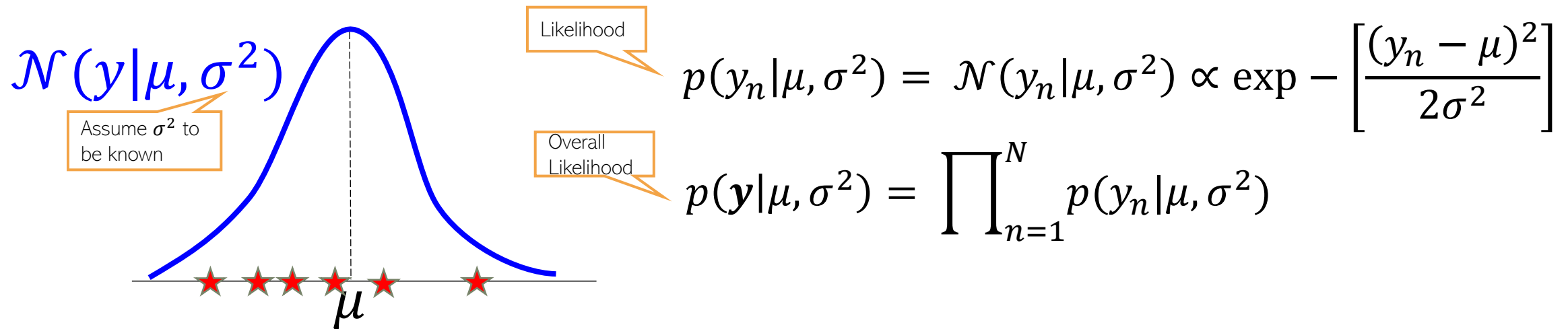
- Many probabilistic ML models are LGMs
- These results are very widely used (PRML Chap. 2 contains a proof)



# Posterior Distribution for Gaussian's Mean

Its MLE/MAP estimation left as an exercise

- Given:  $N$  i.i.d. scalar observations  $\mathbf{y} = \{y_1, y_2, \dots, y_N\}$  assumed drawn from  $\mathcal{N}(y|\mu, \sigma^2)$



- Note: Easy to see that each  $y_n$  drawn from  $\mathcal{N}(y|\mu, \sigma^2)$  is equivalent to the following

Thus  $y_n$  is like a noisy version of  $\mu$  with zero mean Gaussian noise added to it

$$y_n = \mu + \epsilon_n \quad \text{where } \epsilon_n \sim \mathcal{N}(0, \sigma^2)$$

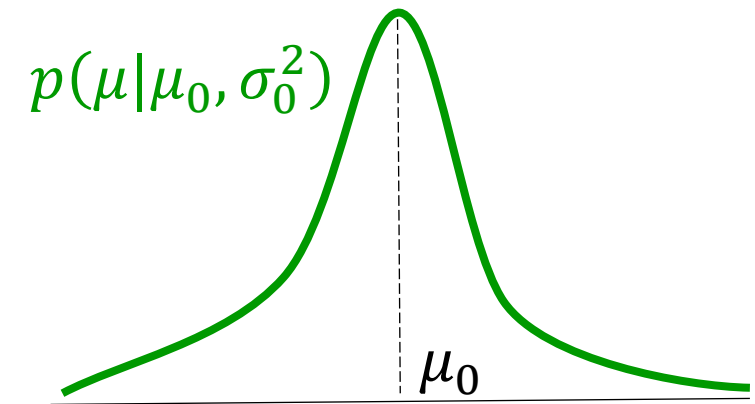
- Let's estimate mean  $\mu$  given  $\mathbf{y}$  using fully Bayesian inference (not point estimation)



# A prior distribution for the mean

- To compute posterior, need a prior over  $\mu$
- Let's choose a Gaussian prior

$$p(\mu|\mu_0, \sigma_0^2) = \mathcal{N}(\mu|\mu_0, \sigma_0^2) \\ \propto \exp \left[ -\frac{(\mu - \mu_0)^2}{2\sigma_0^2} \right]$$



- The prior basically says that a priori we believe  $\mu$  is close to  $\mu_0$
- The prior's variance  $\sigma_0^2$  denotes how certain we are about our belief
- We will assume that the prior's hyperparameters  $(\mu_0, \sigma_0^2)$  are known
- Since  $\sigma^2$  in the likelihood  $\mathcal{N}(y|\mu, \sigma^2)$  is known, Gaussian prior  $\mathcal{N}(\mu|\mu_0, \sigma_0^2)$  on  $\mu$  is also conjugate to the likelihood (thus posterior of  $\mu$  will also be Gaussian)

# The posterior distribution for the mean

- The posterior distribution for the unknown mean parameter  $\mu$

On conditioning side, skipping all fixed params and hyperparams from the notation

$$p(\mu|\mathbf{y}) = \frac{p(\mathbf{y}|\mu)p(\mu)}{p(\mathbf{y})} \propto \prod_{n=1}^N \exp \left[ -\frac{(y_n - \mu)^2}{2\sigma^2} \right] \exp \left[ -\frac{(\mu - \mu_0)^2}{2\sigma_0^2} \right]$$

- Easy to see that the above will be prop. to **exp of a quadratic function** of  $\mu$ . Simplifying:

$$p(\mu|\mathbf{y}) \propto \exp \left[ -\frac{(\mu - \mu_N)^2}{2\sigma_N^2} \right]$$

Gaussian posterior's precision is the sum of the prior's precision and sum of the noise precisions of all the observations

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

Contribution from the prior

Contribution from the data

Also the MLE solution for  $\mu$

Gaussian posterior's mean is a convex combination of prior's mean and the MLE solution

$$\mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \bar{y} \quad (\text{where } \bar{y} = \frac{\sum_{n=1}^N y_n}{N})$$

- What happens to the posterior as  $N$  (number of observations) grows very large?

- Data (likelihood part) overwhelms the prior
- Posterior's variance  $\sigma_N^2$  will approximately be  $\sigma^2/N$  (and goes to 0 as  $N \rightarrow \infty$ )
- The posterior's mean  $\mu_N$  approaches  $\bar{y}$  (which is also the MLE solution)

Meaning, we become very-very certain about the estimate of  $\mu$

# The Predictive Distribution

- If given a point estimate  $\hat{\mu}$ , the plug-in predictive distribution for a test  $y_*$  would be

This is an approximation of the true PPD  $p(y_*|\mathbf{y})$

The best point estimate

$$p(y_*|\hat{\mu}, \sigma^2) = \mathcal{N}(y_*|\hat{\mu}, \sigma^2)$$

- On the other hand, the posterior predictive distribution of  $y_*$  would be

$$\begin{aligned} p(y_*|\mathbf{y}) &= \int p(y_*|\mu, \sigma^2)p(\mu|\mathbf{y})d\mu \\ &= \int \mathcal{N}(y_*|\mu, \sigma^2)\mathcal{N}(\mu|\mu_N, \sigma_N^2)d\mu \\ &= \mathcal{N}(y_*|\mu_N, \sigma^2 + \sigma_N^2) \end{aligned}$$

This "extra" variance  $\sigma_N^2$  in PPD is due to the averaging over the posterior's uncertainty

If **conditional** is Gaussian then **marginal** is also Gaussian

**A useful fact:** When we have conjugacy, the posterior predictive distribution also has a closed form (will see this result more formally when talking about exponential family distributions)



- For an alternative way to get the above result, note that, for test data

$$y_* = \mu + \epsilon \quad \mu \sim \mathcal{N}(\mu_N, \sigma_N^2) \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Using the **posterior** of  $\mu$  since we are at test stage now

$$\Rightarrow p(y_*|\mathbf{y}) = \mathcal{N}(y_*|\mu_N, \sigma^2 + \sigma_N^2)$$

Since both  $\mu$  and  $\epsilon$  are Gaussian r.v., and are independent,  $y_*$  also has a Gaussian posterior predictive, and the respective means and variances of  $\mu$  and  $\epsilon$  get added up

PRML [Bis 06], 2.115, and also mentioned in prob-stats refresher slides



# Gaussian Observation Model: Some Other Facts

- MLE/MAP for  $\mu, \sigma^2$  (or both) is straightforward in Gaussian observation models.
- Posterior also straightforward in most situations for such models
  - (As we saw) computing posterior of  $\mu$  is easy (using Gaussian prior) if variance  $\sigma^2$  is known
  - Likewise, computing posterior of  $\sigma^2$  is easy (using **gamma prior** on  $\sigma^2$ ) if mean  $\mu$  is known
- If  $\mu, \sigma^2$  both are unknown, posterior computation requires computing  $p(\mu, \sigma^2 | \mathbf{y})$ 
  - Computing joint posterior  $p(\mu, \sigma^2 | \mathbf{y})$  exactly requires a jointly conjugate prior  $p(\mu, \sigma^2)$
  - “**Gaussian-gamma**” (“Normal-gamma”) is such a conjugate prior – a product of normal and gamma
  - Note: Computing joint posteriors exactly is possible only in rare cases such this one
- If each observation  $\mathbf{y}_n \in \mathbb{R}^D$ , can assume a likelihood/observation model  $\mathcal{N}(\mathbf{y} | \mu, \Sigma)$ 
  - Need to estimate a **vector-valued** mean  $\mu \in \mathbb{R}^D$ . Can use a **multivariate Gaussian prior**
  - Need to estimate a  $D \times D$  positive definite covariance **matrix**  $\Sigma$ . Can use a **Wishart prior**
  - If  $\mu, \Sigma$  both are unknown, can use **Normal-Wishart** as a conjugate prior

