

ECE595 / STAT598: Machine Learning I

Lecture 11 Maximum-Likelihood Estimation

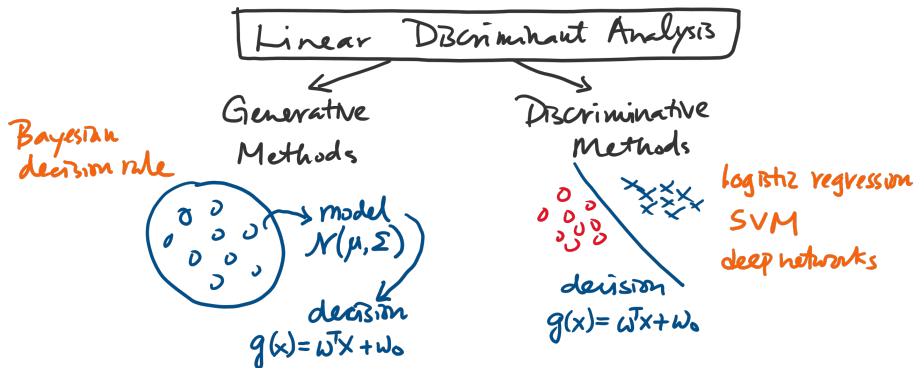
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Overview



- In linear discriminant analysis (LDA), there are generally two types of approaches
- **Generative approach:** Estimate model, then define the classifier
- **Discriminative approach:** Directly define the classifier

Outline

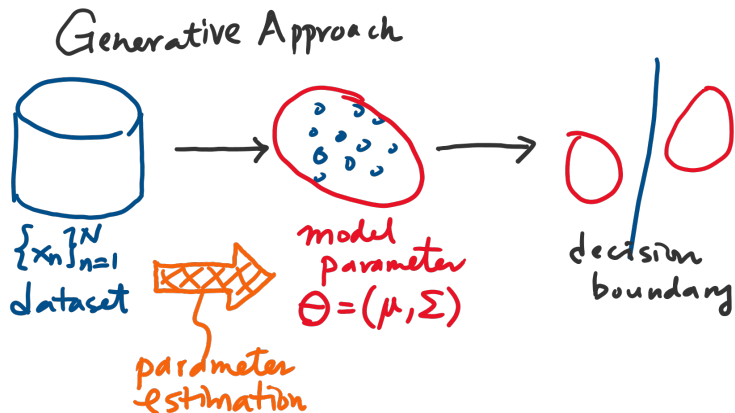
Generative Approaches

- Lecture 9 Bayesian Decision Rules
- Lecture 10 Evaluating Performance
- **Lecture 11 Parameter Estimation**
- Lecture 12 Bayesian Prior
- Lecture 13 Connecting Bayesian and Linear Regression

Today's Lecture

- Basic Principles
 - Likelihood Function
 - Maximum Likelihood Estimate
 - 1D Illustration
 - Gaussian Distributions
- Examples
 - Non-Gaussian Distributions
 - Biased and Unbiased Estimators
 - From MLE to MAP

What is Parameter Estimation?



- The goal of parameter estimation is to determine $\Theta = (\mu, \Sigma)$ from dataset
- This is *the step* where you use data

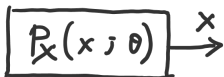
MLE and MAP

There are two typical ways of estimating parameters.

The Generative Process



Bayesian
(MAP estimation)



Frequentist
(ML estimation)

- Maximum-likelihood estimation (MLE): θ is deterministic.
- Maximum-a-posteriori estimation (MAP): θ is random and has a prior distribution.

Maximum Likelihood Estimation

Given the dataset $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$, how to estimate the model parameters?

- We are going to use Gaussian as an illustration.
- Denote θ as the model parameter.
- In Gaussian

$$\theta = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$$

- The likelihood for one data point \mathbf{x}_n is

$$p(\mathbf{x}_n | \overbrace{\theta}^{\{ \boldsymbol{\mu}, \boldsymbol{\Sigma} \}}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\}$$

- θ is a deterministic quantity, not a random variable.
- θ does not have a distribution.
- θ is fixed but unknown.

Likelihood for the Entire Dataset

- Likelihood for the entire dataset $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ is

$$\begin{aligned} p(\mathcal{D} | \boldsymbol{\theta}) &= \prod_{n=1}^N \left\{ \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\} \right\} \\ &= \left(\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \right)^N \exp \left\{ \sum_{n=1}^N -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\} \end{aligned}$$

- The Negative Log-Likelihood is

$$\begin{aligned} -\log p(\mathcal{D} | \boldsymbol{\theta}) &= \frac{N}{2} \log |\boldsymbol{\Sigma}| + \frac{N}{2} \log (2\pi)^d \\ &\quad + \sum_{n=1}^N \left\{ \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\}. \end{aligned}$$

Maximum Likelihood Estimation

- Goal: Find θ that maximizes the likelihood:

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} p(\mathcal{D} | \theta) \\ &= \operatorname{argmax}_{\theta} \prod_{n=1}^N \left\{ \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\} \right\} \\ &= \operatorname{argmin}_{\theta} -\log(\dots) \\ &= \operatorname{argmin}_{\theta} \frac{N}{2} \log |\Sigma| + \frac{N}{2} \log(2\pi)^d \\ &\quad + \sum_{n=1}^N \left\{ \frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\}.\end{aligned}$$

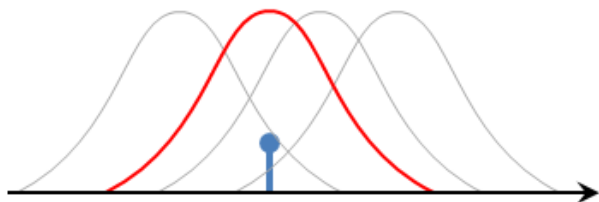
- This optimization is called the **maximum likelihood estimation** (MLE).

Illustrating MLE when $N = 1$. Known σ .

When $N = 1$: The MLE solution is

$$\begin{aligned}\hat{\mu} &= \operatorname{argmax}_{\mu} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_1 - \mu)^2}{2\sigma^2}\right\} \\ &= \operatorname{argmin}_{\mu} (x_1 - \mu)^2 = x_1.\end{aligned}$$

- Which μ will give you the best Gaussian?
- When $\mu = x_1$, the probability of obtaining x_1 is the highest.

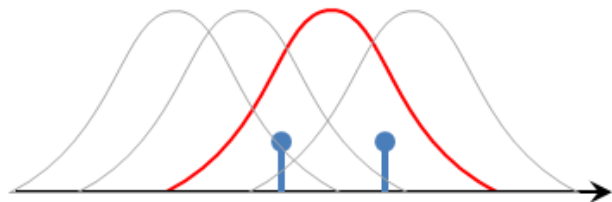


Illustrating MLE when $N = 2$. Known σ .

When $N = 2$: The MLE solution is

$$\begin{aligned}\hat{\mu} &= \operatorname{argmax}_{\mu} \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^2 \exp \left\{ -\frac{(x_1 - \mu)^2 + (x_2 - \mu)^2}{2\sigma^2} \right\} \\ &= \operatorname{argmin}_{\mu} (x_1 - \mu)^2 + (x_2 - \mu)^2 = \frac{x_1 + x_2}{2}.\end{aligned}$$

- Which μ will give you the best Gaussian?
- When $\mu = (x_1 + x_2)/2$, the prob. of obtaining x_1 and x_2 is highest.

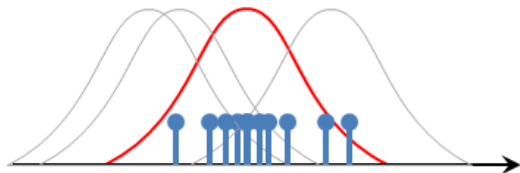


Illustrating MLE when $N =$ arbitrary integer

The MLE solution is

$$\begin{aligned}\hat{\mu} &= \operatorname{argmax}_{\mu} \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^2 \exp \left\{ - \sum_{n=1}^N \frac{(x_n - \mu)^2}{2\sigma^2} \right\} \\ &= \operatorname{argmin}_{\mu} \sum_{n=1}^N (x_n - \mu)^2 = \frac{1}{N} \sum_{n=1}^N x_n.\end{aligned}$$

- Which μ will give you the best Gaussian?
- When $\mu = \frac{1}{N} \sum_{n=1}^N x_n$, the prob. of obtaining $\{x_n\}$ is highest.



Estimation in High-dimension

- Assume Σ is known and fixed.
- Thus, $\theta = \mu$. Estimate μ

$$\begin{aligned}\hat{\mu} &= \operatorname{argmin}_{\mu} \cancel{\frac{N}{2} \log |\Sigma|} + \cancel{\frac{N}{2} \log (2\pi)^d} \\ &\quad + \sum_{n=1}^N \left\{ \frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\} \\ &= \operatorname{argmin}_{\mu} \sum_{n=1}^N \left\{ (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\}\end{aligned}$$

- Take derivative, setting to zero:

$$\nabla_{\mu} \left\{ \sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\} = 2 \sum_{n=1}^N \Sigma^{-1} (\mathbf{x}_n - \mu) = \mathbf{0}.$$

Estimation in High-dimension

- Let us do some algebra

$$\sum_{n=1}^N \Sigma^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) = \mathbf{0} \implies \sum_{n=1}^N \mathbf{x}_n = \sum_{n=1}^N \boldsymbol{\mu}$$

- Then we can show that the MLE solution is

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$

- This is just the **empirical average** of the entire dataset!
- You can show that if $\mathbb{E}[\mathbf{x}_n] = \boldsymbol{\mu}$ for all n , then

$$\mathbb{E}[\hat{\boldsymbol{\mu}}] = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[\mathbf{x}_n] = \boldsymbol{\mu}.$$

- We say that $\hat{\boldsymbol{\mu}}$ is a **unbiased estimator** of $\boldsymbol{\mu}$ since $\mathbb{E}[\hat{\boldsymbol{\mu}}] = \boldsymbol{\mu}$.

When both μ and Σ are Unknown

What will be the MLE when both μ and Σ are unknown?

$$\begin{aligned}(\hat{\mu}, \hat{\Sigma}) = \operatorname{argmin}_{\mu, \Sigma} & \frac{N}{2} \log |\Sigma| + \frac{N}{2} \log(2\pi)^d \\ & + \underbrace{\sum_{n=1}^N \left\{ \frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\}}_{\varphi(\mu, \Sigma)}.\end{aligned}$$

You need to take derivative with respect to μ and Σ , and solve

$$\begin{aligned}\nabla_{\mu} \varphi(\mu, \Sigma) &= \mathbf{0} \\ \nabla_{\Sigma} \varphi(\mu, \Sigma) &= \mathbf{0}\end{aligned}$$

With some (tedious) matrix calculus, we can show that

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \quad \text{and} \quad \hat{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \hat{\mu})(\mathbf{x}_n - \hat{\mu})^T.$$

Exercise: Prove this result when \mathbf{x}_n is a 1D scalar.

Outline

Generative Approaches

- Lecture 9 Bayesian Decision Rules
- Lecture 10 Evaluating Performance
- **Lecture 11 Parameter Estimation**
- Lecture 12 Bayesian Prior
- Lecture 13 Connecting Bayesian and Linear Regression

Today's Lecture

- Basic Principles
 - Likelihood Function
 - Maximum Likelihood Estimate
 - 1D Illustration
 - Gaussian Distributions
- Examples
 - Non-Gaussian Distributions
 - Biased and Unbiased Estimators
 - From MLE to MAP

MLE does not need to be Gaussian

Suppose that $x_n \sim \text{Bernoulli}(\theta)$. Then,

$$p(x_n|\theta) = \begin{cases} \theta, & \text{if } x_n = 1 \\ 1 - \theta, & \text{if } x_n = 0. \end{cases}$$

We can write the likelihood function as

$$\begin{aligned} p(\mathcal{D}|\theta) &= \prod_{n=1}^N \theta^{x_n} (1 - \theta)^{1-x_n} \\ &= \theta^{\sum_{n=1}^N x_n} (1 - \theta)^{\sum_{n=1}^N (1-x_n)} \end{aligned}$$

and so the negative-log likelihood is

$$\begin{aligned} -\log p(\mathcal{D}|\theta) &= -\log \left\{ \theta^{\sum_{n=1}^N x_n} (1 - \theta)^{\sum_{n=1}^N (1-x_n)} \right\} \\ &= - \left(\sum_{n=1}^N x_n \right) \log \theta - \left(\sum_{n=1}^N (1 - x_n) \right) \log(1 - \theta) \end{aligned}$$

Bernoulli MLE

Taking the derivative and setting to zero:

$$\begin{aligned}\frac{d}{d\theta} \left\{ -\log p(\mathcal{D}|\theta) \right\} &= \frac{d}{d\theta} \left\{ - \left(\sum_{n=1}^N x_n \right) \log \theta - \left(\sum_{n=1}^N (1 - x_n) \right) \log(1 - \theta) \right\} \\ &= - \left(\sum_{n=1}^N x_n \right) \frac{1}{\theta} + \left(\sum_{n=1}^N (1 - x_n) \right) \frac{1}{1 - \theta}.\end{aligned}$$

Setting to zero, we have

$$\begin{aligned}\left(\sum_{n=1}^N x_n \right) \frac{1}{\theta} &= \left(\sum_{n=1}^N (1 - x_n) \right) \frac{1}{1 - \theta} \\ \left(\sum_{n=1}^N x_n \right) (1 - \theta) &= \left(\sum_{n=1}^N (1 - x_n) \right) \theta \\ \left(\sum_{n=1}^N x_n \right) - \left(\sum_{n=1}^N x_n \right) \theta &= N\theta - \left(\sum_{n=1}^N x_n \right) \theta\end{aligned}$$

Therefore,

$$\theta = \frac{1}{N} \sum_{n=1}^N x_n.$$

Unbias and Consistent Estimator

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \quad \text{and} \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \hat{\boldsymbol{\mu}})(\mathbf{x}_n - \hat{\boldsymbol{\mu}})^T.$$

- $\hat{\boldsymbol{\mu}}$ is the **empirical average** (or the sample mean)
- $\hat{\boldsymbol{\Sigma}}$ is the **empirical covariance** (or the sample covariance)
- $\mathbb{E}[\hat{\boldsymbol{\mu}}] = \boldsymbol{\mu}$
 - $\hat{\boldsymbol{\mu}}$ is an **unbiased** estimate of $\boldsymbol{\mu}$: $\mathbb{E}[\hat{\boldsymbol{\mu}}] = \boldsymbol{\mu}$ for all N
 - $\hat{\boldsymbol{\mu}}$ is a consistent estimate of $\boldsymbol{\mu}$: As $N \rightarrow \infty$, $\hat{\boldsymbol{\mu}} \xrightarrow{p} \boldsymbol{\mu}$
- $\mathbb{E}[\hat{\boldsymbol{\Sigma}}] = \frac{N-1}{N} \boldsymbol{\Sigma}$
 - $\hat{\boldsymbol{\Sigma}}$ is a **biased** estimate of $\boldsymbol{\Sigma}$: $\mathbb{E}[\hat{\boldsymbol{\Sigma}}] \neq \boldsymbol{\Sigma}$
 - $\hat{\boldsymbol{\Sigma}}$ is a consistent estimate of $\boldsymbol{\Sigma}$: As $N \rightarrow \infty$, $\hat{\boldsymbol{\Sigma}} \xrightarrow{p} \boldsymbol{\Sigma}$
- You can make $\boldsymbol{\Sigma}$ unbiased by defining

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \quad \text{and} \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{N-1} \sum_{n=1}^N (\mathbf{x}_n - \hat{\boldsymbol{\mu}})(\mathbf{x}_n - \hat{\boldsymbol{\mu}})^T.$$

Unbiased and Consistent Estimator

Where does $(N - 1)/N$ come from? Here is a 1D explanation. Assume $\mu = 0$.

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N x_n, \quad \text{and} \quad \hat{\sigma}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \hat{\mu})^2$$

Taking expectation yields

$$\begin{aligned} \mathbb{E}[\hat{\sigma}^2] &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}(x_n - \hat{\mu})^2 \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{E}[x_n^2] - 2\mathbb{E}[\hat{\mu}x_n] + \mathbb{E}[\hat{\mu}^2] \\ &= \frac{1}{N} \sum_{n=1}^N \left\{ \sigma^2 - 2\mathbb{E} \left[\frac{1}{N} \sum_{j=1}^N x_j x_n \right] + \mathbb{E} \left[\left(\frac{1}{N} \sum_{n=1}^N x_n \right)^2 \right] \right\} \end{aligned}$$

Unbiased and Consistent Estimator

$$\begin{aligned} \mathbb{E}[\hat{\sigma}^2] &= \frac{1}{N} \sum_{n=1}^N \left\{ \sigma^2 - \underbrace{2\mathbb{E}\left[\frac{1}{N} \sum_{j=1}^N x_j x_n\right]}_{\substack{= \frac{2}{N}(\mathbb{E}[x_1 x_n + \dots + x_N x_n]) \\ = \frac{2}{N}(0 + \dots + \sigma^2 + \dots + 0) \\ = \frac{2\sigma^2}{N}}} + \underbrace{\mathbb{E}\left[\left(\frac{1}{N} \sum_{n=1}^N x_n\right)^2\right]}_{\substack{= \frac{1}{N^2} \sum_{n=1}^N \mathbb{E}[x_n^2] + \sum_{j \neq n} \mathbb{E}[x_j x_n] \\ = \frac{1}{N^2} N\sigma^2 + 0 \\ = \frac{\sigma^2}{N}}} \right\} \\ &= \frac{1}{N} \sum_{n=1}^N \left\{ \sigma^2 - \frac{2}{N} \sigma^2 + \frac{1}{N} \sigma^2 \right\} \\ &= \frac{N-1}{N} \sigma^2. \end{aligned}$$

From ML to Decision Boundary

- If you have a training set \mathcal{D} , ...
- Partition it according to the labels: $\mathcal{D}^{(1)}$ and $\mathcal{D}^{(2)}$
- Pick a model, e.g., Gaussian
- For each class, estimate the model parameter
 - $\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1$ for $\mathcal{D}^{(1)}$
 - $\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2$ for $\mathcal{D}^{(2)}$
- Construct a discriminant function

$$g(\mathbf{x}) = \mathbf{x}^T (\mathbf{W}_1 - \mathbf{W}_2) \mathbf{x} + (\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x} + (w_{10} - w_{20})$$

- See Lecture Bayesian Decision Rule 1 for formula
- Define the hypothesis function

$$h(\mathbf{x}) = \begin{cases} 1, & \text{if } g(\mathbf{x}) > 0, \\ 0, & \text{if } g(\mathbf{x}) < 0. \end{cases}$$

How well do you do?

- You have a dataset $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$.
- You estimate $\boldsymbol{\mu}$ by solving the maximum-likelihood problem.
- You get an estimate

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n$$

- How good is your estimate?
- Hoeffding inequality. In the 1D case:

$$\mathbb{P}[|\hat{\mu} - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}.$$

- As $N \rightarrow \infty$, $\hat{\mu} \rightarrow \mu$ with very high probability.
- Overall performance:
 - Is $\hat{\mu} \approx \mu$?
 - Is μ giving a good classifier?

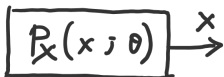
MLE and MAP

There are two typical ways of estimating parameters.

The Generative Process



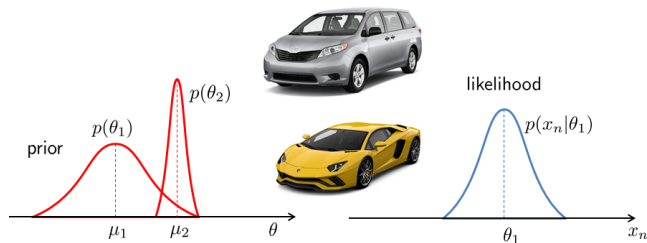
Bayesian
(MAP estimation)



Frequentist
(ML estimation)

- Maximum-likelihood estimation (MLE): θ is deterministic.
- Maximum-a-posteriori estimation (MAP): θ is random and has a prior distribution.

From MLE to MAP



- Likelihood:

$$p(x_n|\theta_1) = \mathcal{N}(x_n|\theta_1, \sigma_1^2), \quad \text{and} \quad p(x_n|\theta_2) = \mathcal{N}(x_n|\theta_2, \sigma_2^2).$$

- Maximum-likelihood: You know nothing about θ_1 and θ_2 . So you need to take measurements to estimate θ_1 and θ_2 .
- Maximum-a-Posteriori: You know something about θ_1 and θ_2 .
- Prior

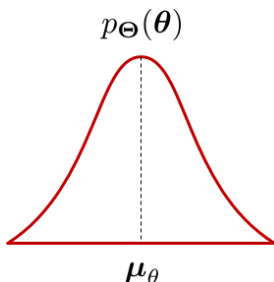
$$p(\theta_1) = \mathcal{N}(\mu_1|\gamma_1^2), \quad \text{and} \quad p(\theta_2) = \mathcal{N}(\mu_2|\gamma_2^2).$$

From MLE to MAP

- In MLE, the parameter θ is **deterministic**.
- What if we assume θ has a distribution?
- This makes θ **probabilistic**.
- So make Θ as a random variable, and θ a state of Θ .
- Distribution of Θ :

$$p_{\Theta}(\theta)$$

- $p_{\Theta}(\theta)$ is the distribution of the parameter Θ .
- Θ has its own mean and own variance.



Maximum-a-Posteriori

By Bayes Theorem again:

$$p_{\Theta|\mathbf{X}}(\theta|\mathbf{x}_n) = \frac{p_{\mathbf{X}|\Theta}(\mathbf{x}_n|\theta)p_{\Theta}(\theta)}{p_{\mathbf{X}}(\mathbf{x}_n)}.$$

- To maximize the posterior distribution

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} p_{\Theta|\mathbf{X}}(\theta|\mathcal{D}) \\ &= \operatorname{argmax}_{\theta} \prod_{n=1}^N p_{\Theta|\mathbf{X}}(\theta|\mathbf{x}_n) \\ &= \operatorname{argmax}_{\theta} \prod_{n=1}^N \frac{p_{\mathbf{X}|\Theta}(\mathbf{x}_n|\theta)p_{\Theta}(\theta)}{p_{\mathbf{X}}(\mathbf{x}_n)} \\ &= \operatorname{argmin}_{\theta} - \sum_{n=1}^N \left\{ \log p_{\mathbf{X}|\Theta}(\mathbf{x}_n|\theta) + \log p_{\Theta}(\theta) \right\}\end{aligned}$$

Reading List

Maximum Likelihood Estimation

- Duda-Hart-Stork, Pattern Classification, Chapter 3.2
- Iowa State EE 527 http://www.ece.iastate.edu/~namrata/EE527_Spring08/15.pdf
- Purdue ECE 645, Lecture 18-20
<https://engineering.purdue.edu/ChanGroup/ECE645.html>
- UCSD ECE 271A, Lecture 6
<http://www.svcl.ucsd.edu/courses/ece271A/ece271A.htm>
- Univ. Orleans. https://www.univ-orleans.fr/deg/masters/ESA/CH/Chapter2_MLE.pdf