

Chapter 8

Statistics for Machine Learning

8.1 Representing Data

- Numeric encoding of data
- 3 types
 - Categorical, Ordinal, Numerical

- 8.1.1 Categorical (or nominal)
 - No intrinsic ordering
 - Example
 - Jobs: soldier, sailor, tinker, spy
 - Representation
 - Integer encoding: 1, 2, 3, 4
 - 1-of- m encoding: (1,0,0,0), (0,1,0,0), (0,0,1,0), (0,0,0,1)

- 8.1.2 Ordinal
 - Discrete
 - Intrinsic ordering or ranking
 - Example
 - cold, cool, warm, hot
 - Representation
 - To preserve the ordering
 - $-1, 0, 1, 2$

- 8.1.3 Numerical
 - Values that are real numbers
 - Example
 - Temperature
 - Salary

8.2 Distributions

- Distributions over discrete variables
 - $\text{dom}(x)$ is discrete
 - Example: dice
- Distributions over continuous variables
 - $\text{dom}(x)$ is an interval or a region

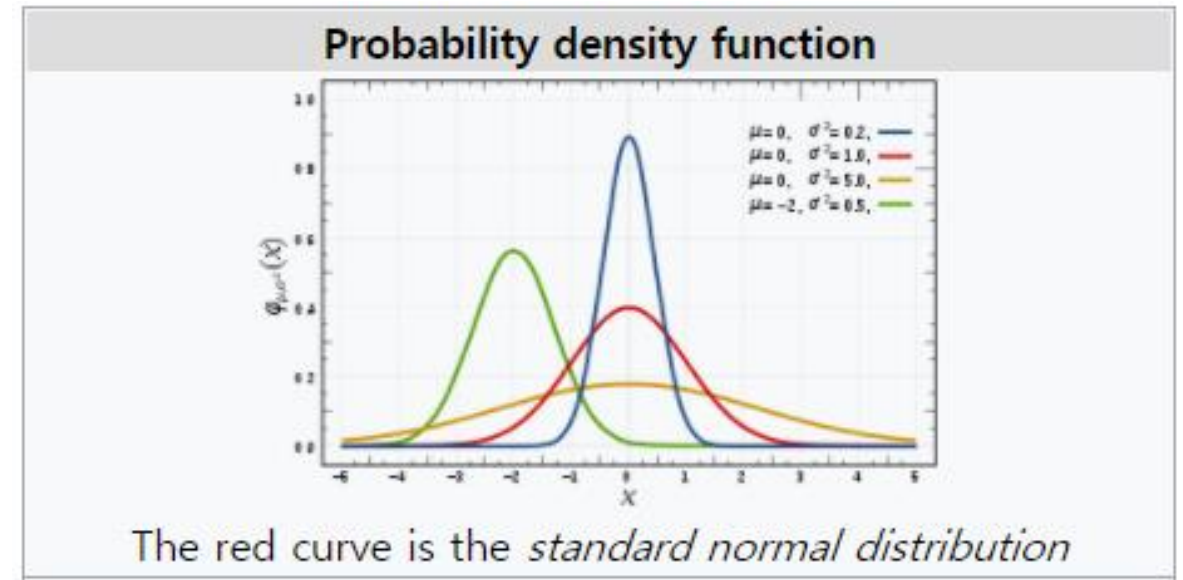
- Definition 8.1 (Probability Density Functions)

- Probability density function $p(x)$

$$p(x) \geq 0, \quad \int_{-\infty}^{\infty} p(x)dx = 1$$

- Probability

$$p(a \leq x \leq b) = \int_a^b p(x)dx$$



- Definition 8.2 (Averages and Expectation)

- Average (or expectation)

- $\langle f(x) \rangle_{p(x)}$

- $\mathbb{E}(f(x))$

Discrete case

$$\langle f(x) \rangle_{p(x)} = \sum_x f(x)p(x)$$

Continuous case

$$\langle f(x) \rangle_{p(x)} = \int_{-\infty}^{\infty} f(x)p(x)dx$$

- Result 8.1 (Change of variables)

1-dimensional x

$$y = f(x)$$

$$p(x) = p(y) \frac{dy}{dx}$$

n -dimensional \mathbf{x}

$$\mathbf{y} = f(\mathbf{x})$$

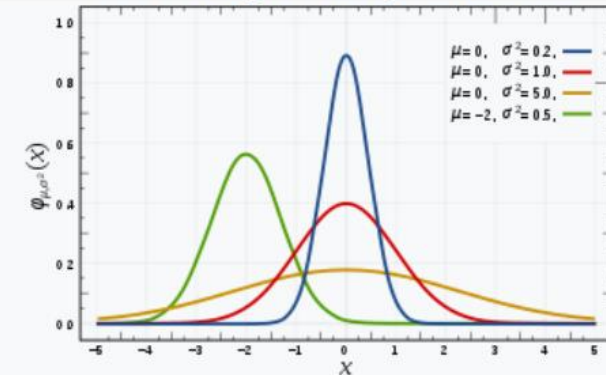
$$p(\mathbf{x}) = p(\mathbf{y}) \left| \det \left(\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \right) \right|$$

- Definition 8.4 (Cumulative Distribution Function)

$$\text{cdf}(t) = p(x < t) = \langle \mathbb{I}[x < t] \rangle_{p(x)}$$

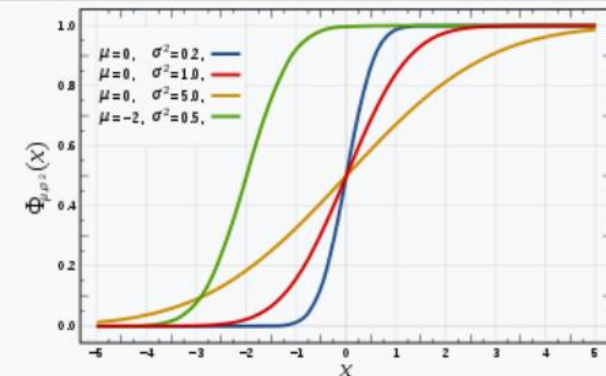
$$\text{cdf}(-\infty) = 0, \text{cdf}(\infty) = 1$$

Probability density function



The red curve is the *standard normal distribution*

Cumulative distribution function



- Definition 8.3 (Moments)

k -th moment

$$\langle x^k \rangle_{p(x)}$$

1st moment = mean

$$\langle x \rangle_{p(x)}$$

- Definition 8.5 (Moment Generating Function)

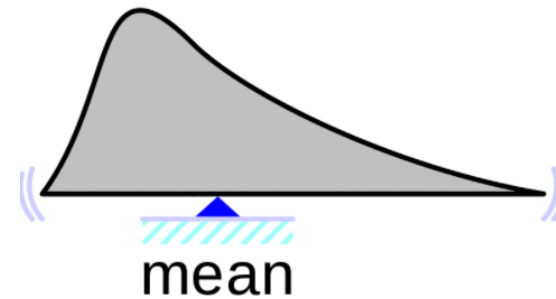
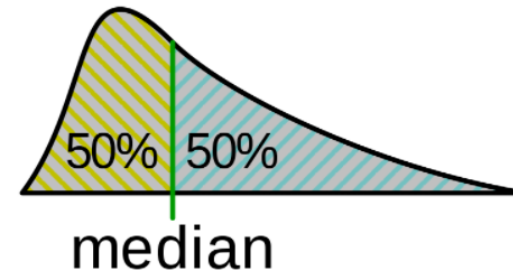
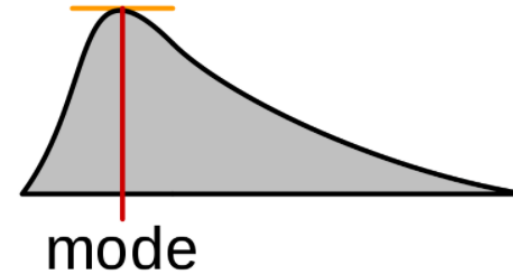
$$g(t) = \langle e^{tx} \rangle_{p(x)} = \int_{-\infty}^{\infty} e^{tx} p(x) dx$$

$$g^{(k)}(0) = \langle x^k \rangle_{p(x)}$$

- Definition 8.6 (Mode)

Highest value

$$x_* = \operatorname{argmax}_x p(x)$$



- Definition 8.7 (Variance and Correlation)

Variance

$$\begin{aligned}\sigma^2 &= \langle (x - \langle x \rangle)^2 \rangle_{p(x)} \\ &= \langle x^2 \rangle - \langle x \rangle^2\end{aligned}$$

Covariance matrix

$$\begin{aligned}\Sigma_{ij} &= \langle (x_i - \mu_i)(x_j - \mu_j) \rangle \\ &= \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle\end{aligned}$$

Correlation matrix

$$\rho_{ij} = \left\langle \frac{x_i - \mu_i}{\sigma_i} \frac{x_j - \mu_j}{\sigma_j} \right\rangle$$

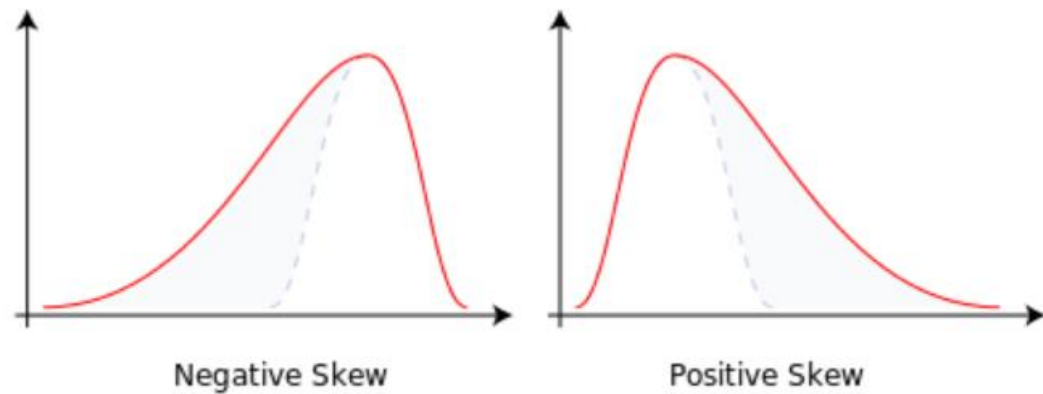
- Definition 8.8 (Skewness and Kurtosis)

Skewness

$$\gamma_1 = \frac{\langle (x - \langle x \rangle)^3 \rangle_{p(x)}}{\sigma^3}$$

Kurtosis

$$\gamma_2 = \frac{\langle (x - \langle x \rangle)^4 \rangle_{p(x)}}{\sigma^4} - 3$$



- Definition 8.10 (Empirical Distribution)

$$\text{dom}(x) = \{x^1, \dots, x^N\}$$

$$p(x) = \frac{1}{N} \sum_{n=1}^N \mathbb{I}[x = x^n]$$

If $x^i = x^j$ for all i and j , then

$$p(x = x^n) = \frac{1}{N}$$

Mean

$$\mu = \frac{1}{N} \sum_{n=1}^N x^n$$

Variance

$$\sigma^2 = \frac{1}{N} \sum_{n=1}^N (x^n - \mu)^2$$

- Definition 8.11 (KL divergence)

Kullback-Leibler divergence between distributions q and p

$$\text{KL}(q|p) = \langle \log q(x) - \log p(x) \rangle_{q(x)}$$

KL divergence is nonnegative

$$\text{KL}(q|p) \geq 0$$

Proof) ?

- Entropy

$$H(p) = -\langle \log p(x) \rangle_{p(x)}$$

$$H(p) = -\text{KL}(p|u) + \text{const.}$$

where u is the uniform distribution

8.3 Classical Distributions

- Definition 8.14 (Bernoulli Distribution)

$$\begin{aligned}\text{dom}(x) &= \{0, 1\} \\ p(x = 1) &= \theta\end{aligned}$$

$$\begin{aligned}\langle x \rangle &= \theta \\ \text{var}(x) &= \theta(1 - \theta)\end{aligned}$$

Example: a coin toss

- Definition 8.15 (Categorical Distribution)

$$\text{dom}(x) = \{1, \dots, C\}$$
$$p(x = c) = \theta_c, \sum_c \theta_c = 1$$

Example: the roll of a dice

- Definition 8.16 (Binomial Distribution)

The probability that in n Bernoulli trials x^1, \dots, x^n
there will be k states 1 observed

$$p(y = k|\theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}$$

where $y = \sum_{i=1}^n \mathbb{I}[x^i = 1]$, the number of 1's

$$\langle y \rangle = n\theta, \text{ var}(y) = n\theta(1 - \theta)$$

Example: coin toss n times

$$\binom{n}{k} = \frac{n!}{k! (n - k)!}$$

- Definition 8.17 (Multinomial Distribution)

n trials of a categorical distribution

Example: n rolls of a dice

- Definition 8.18 (Poisson Distribution)

- The probability of a given number of events occurring in an interval of time

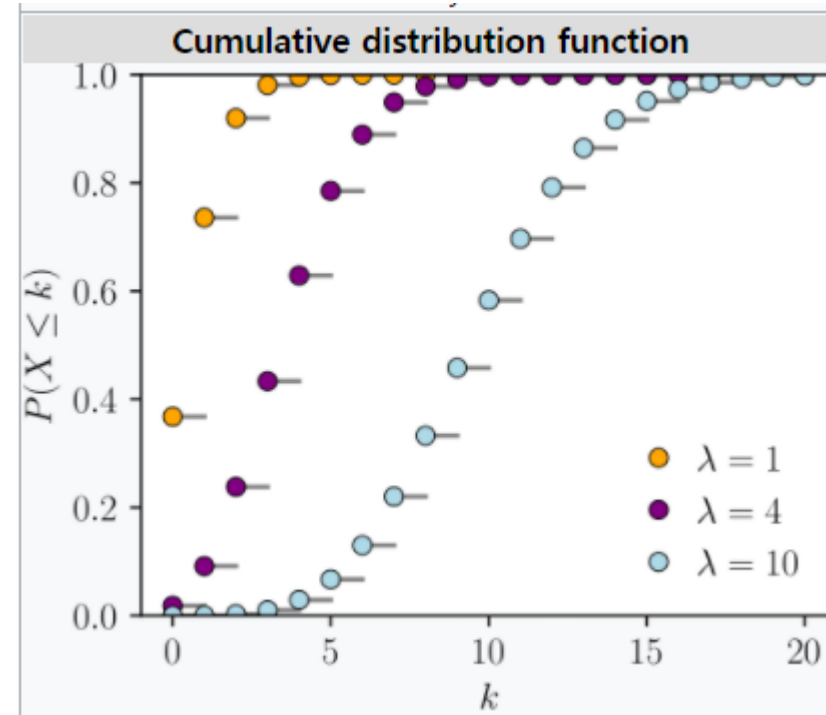
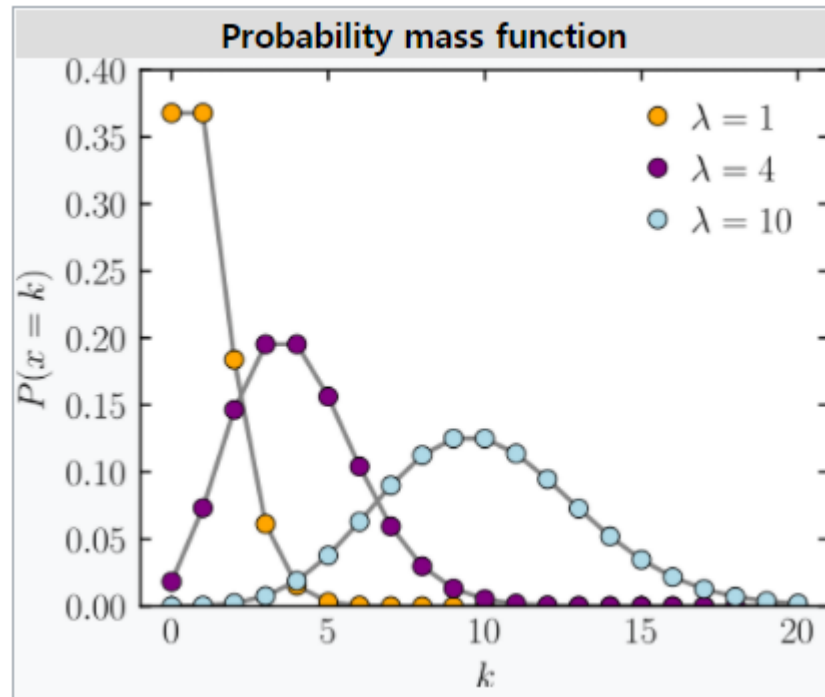
$$p(x = k|\lambda) = \frac{\lambda^k e^{-\lambda}}{k!}$$

$$\begin{array}{l} \langle x \rangle = \lambda, \\ \text{var}(x) = \lambda \end{array}$$

- Example: train

If λ is the expected number of events per unit interval, then the probability of the number of events x within an interval t is

$$p(x = k|\lambda) = \frac{1}{k!} e^{-\lambda t} (\lambda t)^k$$



- Definition 8.19 (Uniform distribution)

$$p(x) = \text{const.}$$

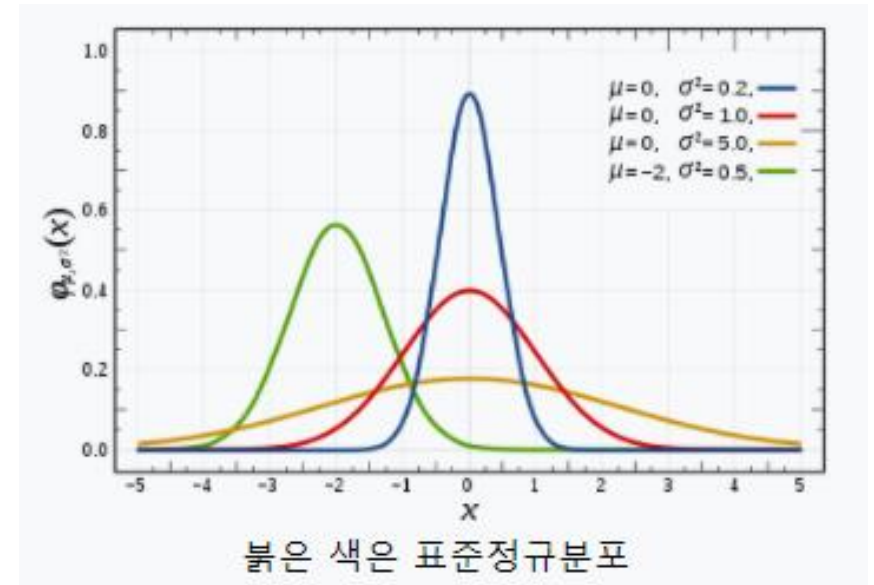
- Definition 8.21 (Gamma Distribution)
- Definition 8.22 (Inverse Gamma distribution)
- Definition 8.23 (Beta Distribution)

- Definition 8.25 (Univariate Gaussian Distribution)

$$\begin{aligned} p(x|\mu, \sigma^2) &= \mathcal{N}(x|\mu, \sigma^2) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \end{aligned}$$

$$\langle x \rangle = \mu$$

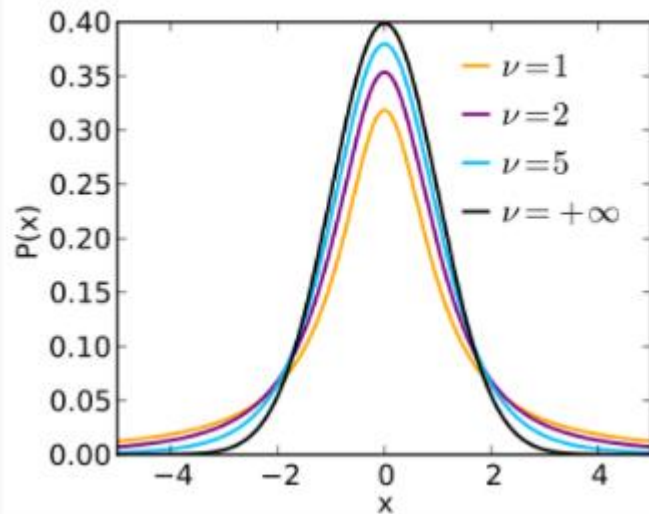
$$\text{var}(x) = \sigma^2$$



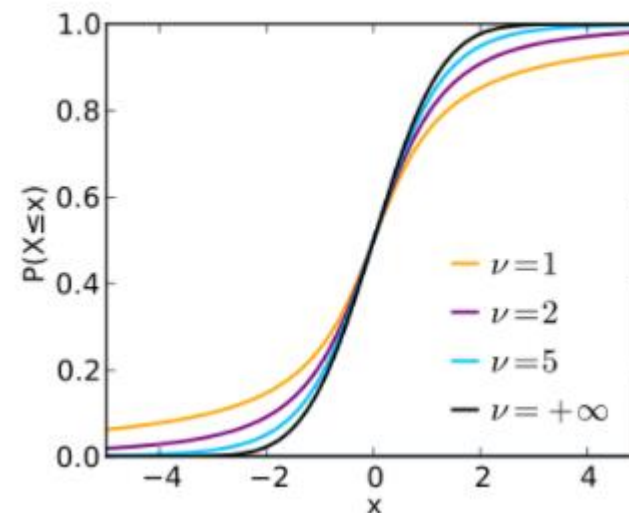
- Definition 8.26 (Student's t-distribution).

$$p(x|\mu, \lambda, \nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)} \left(\frac{\lambda}{\nu\pi}\right)^{\frac{1}{2}} \left(1 + \frac{\lambda(x-\mu)^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

확률 밀도 함수



누적 분포 함수

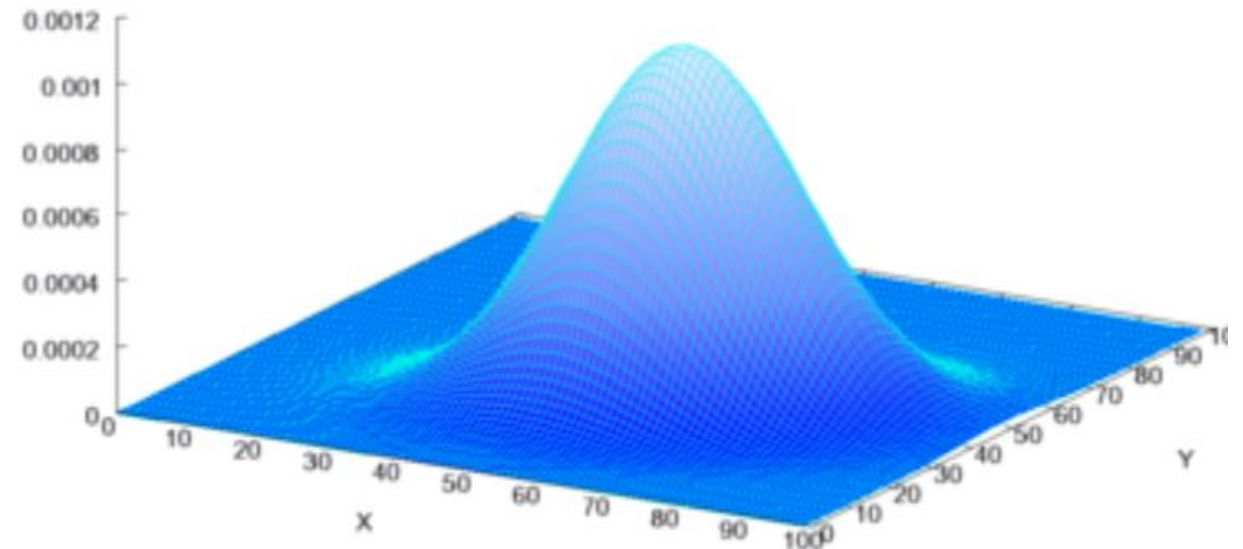


8.4 Multivariate Gaussian

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{\det 2\pi\boldsymbol{\Sigma}}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

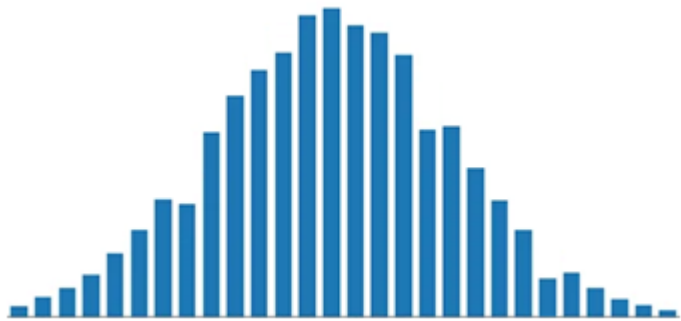
The mean vector $\boldsymbol{\mu}$

The covariance matrix $\boldsymbol{\Sigma}$



8.6 Learning distributions

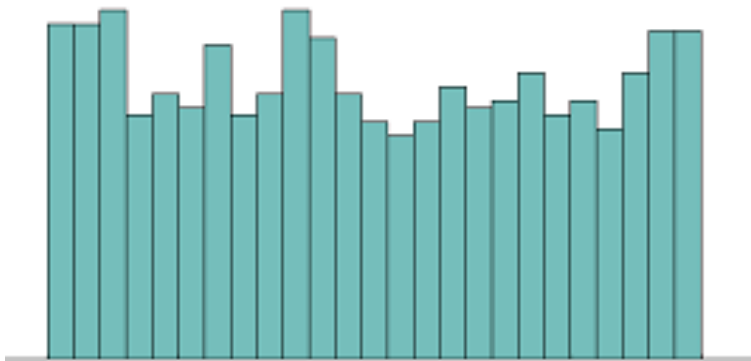
- Learning is Inferring the distribution from data
 - Inferring the distribution $p(x|\theta)$ from $\{x^1, \dots, x^N\}$
 - i.e. determining θ



Assume that the distribution is gaussian

$$\mathcal{N}(x|\mu, \sigma^2)$$

Determine μ and θ from data



Uniform distribution $p(x) = c$
Inferring is to determine c

- Maximum A posteriori
 - $\theta^{MAP} = \operatorname{argmax}_{\theta} p(\theta|\mathcal{X})$
- Maximum Likelihood
 - $\theta^{ML} = \operatorname{argmax}_{\theta} p(\mathcal{X}|\theta)$
- Moment Matching
 - θ is set such that the moment of the distribution matches the empirical moment
- Pseudo Likelihood

- Definition 8.30. Prior, Likelihood and Posterior

$$\underbrace{p(\theta|\mathcal{X})}_{\text{posterior}} = \frac{\underbrace{p(\mathcal{X}|\theta)}_{\text{likelihood}} \underbrace{p(\theta)}_{\text{prior}}}{\underbrace{p(\mathcal{X})}_{\text{evidence}}}$$

$$p(\theta|\mathcal{X}, M) = \frac{p(\mathcal{X}|\theta, M)p(\theta|M)}{p(\mathcal{X}|M)}$$

M is a model, such as gaussian

The most probable a posteriori (MAP) setting is that which maximises the posterior

$$\theta^{MAP} = \operatorname{argmax}_{\theta} p(\theta|\mathcal{X}, M)$$

The maximum likelihood (ML) setting is that which maximises the likelihood

$$\theta^{ML} = \operatorname{argmax}_{\theta} p(\mathcal{X}|\theta, M)$$

A coin is known to be $p(H) = p(T) = 0.5$

Experiment: toss a coin 10 times, 6 heads and 4 tails

Infer $p(H)$

Let $\theta = p(H)$ and $\mathcal{X} = \{H, H, H, H, H, H, T, T, T, T\}$

Model = the binomial distribution

ML

$$p(\mathcal{X}|\theta) = \theta^6(1 - \theta)^4$$

$$\operatorname{argmax}_{\theta} p(\mathcal{X}|\theta) = 0.4$$

MAP

$$p(\theta|\mathcal{X}) = p(\mathcal{X}|\theta)p(\theta)/p(\mathcal{X})$$

$$\sim p(\mathcal{X}|\theta)p(\theta)$$

depends on the distribution of θ