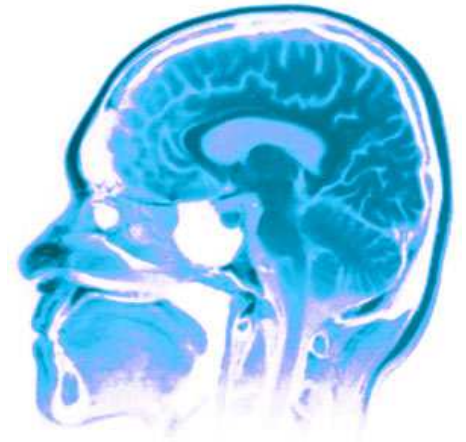# CPSC540

# Probabilistic linear prediction and maximum likelihood

Nando de Freitas

*January, 2013*

*University of British Columbia*

# Outline of the lecture

In this lecture, we formulate the problem of linear prediction using probabilities. We also introduce the maximum likelihood estimate and show that it coincides with the least squares estimate. The goal of the lecture is for you to learn:

- ❑ Multivariate Gaussian distributions
- ❑ How to formulate the likelihood for linear regression
- ❑ Computing the maximum likelihood estimates for linear regression.
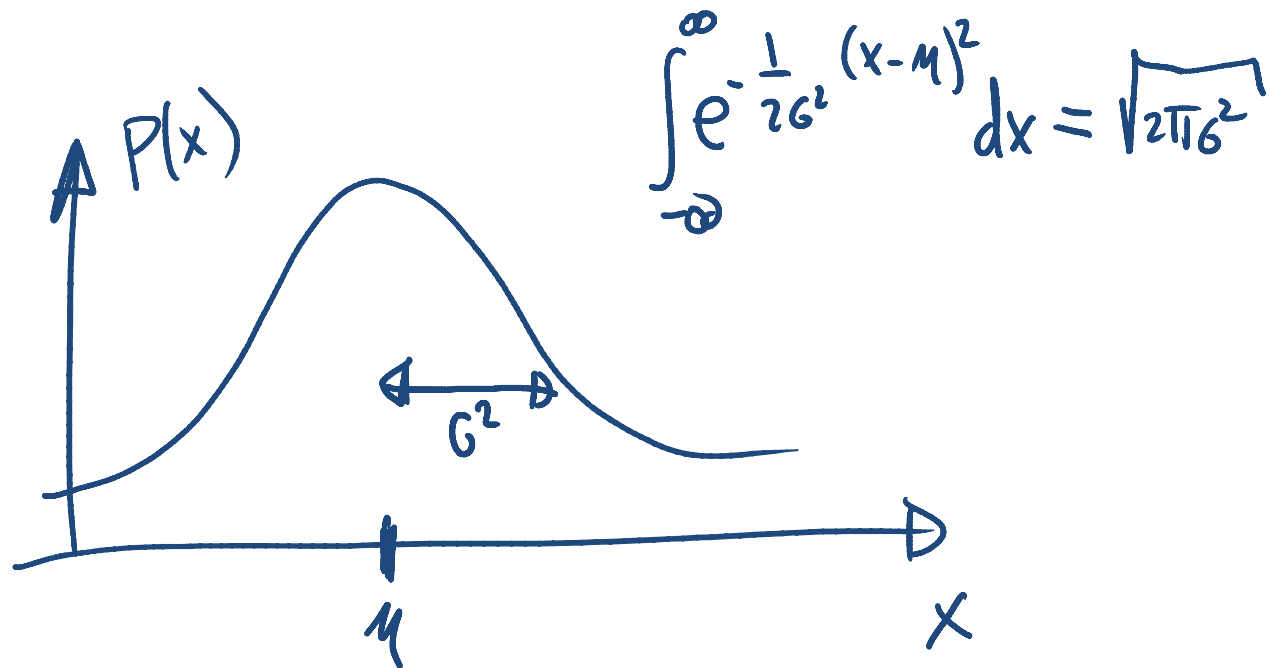- ❑ Understand why maximum likelihood is used.

# Univariate Gaussian distribution

The probability density function (pdf) of a Gaussian distribution is given by
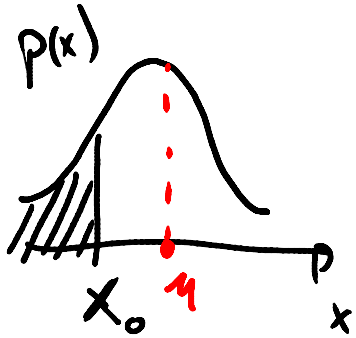
$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}. \qquad x \sim \mathcal{N}(\mu, \sigma^2)$$

where $\mu$ is the mean or center of mass and $\sigma^2$ is the variance.

$$\int_{-\infty}^{\infty} P(x)\,dx = 1$$
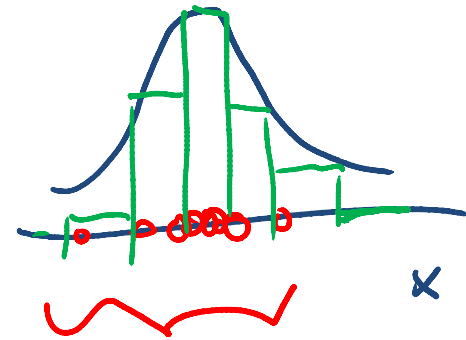
$$\int_{-\infty}^{\infty} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}\,dx = \sqrt{2\pi\sigma^2}$$
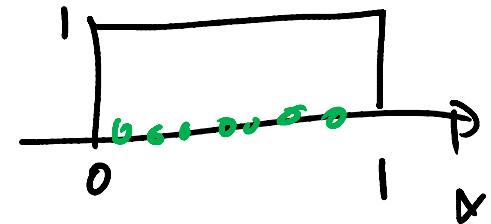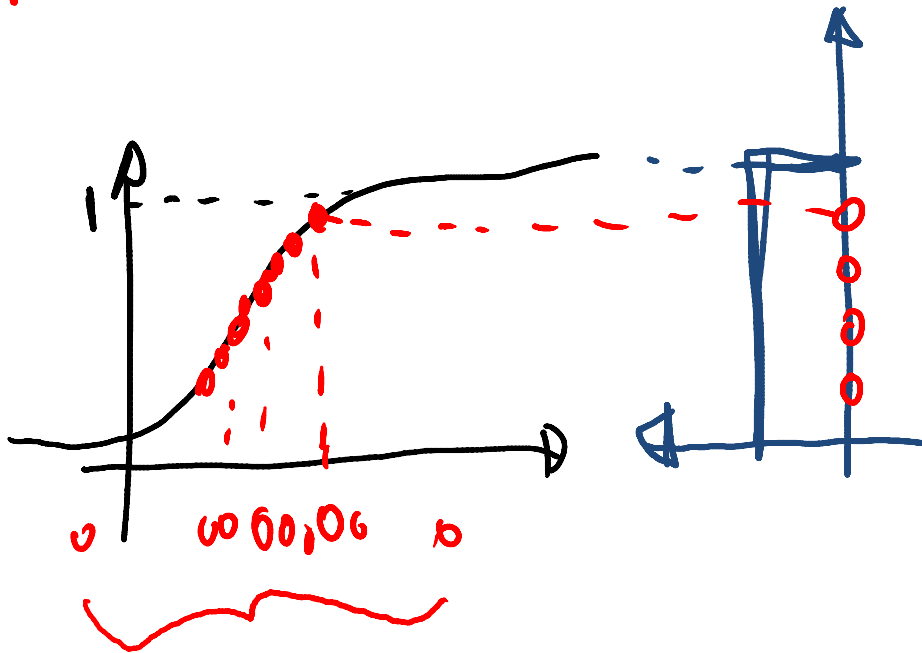
$P(x)$

$\sigma^2$

$\mu$

$x$

# Sampling from a Gaussian distribution



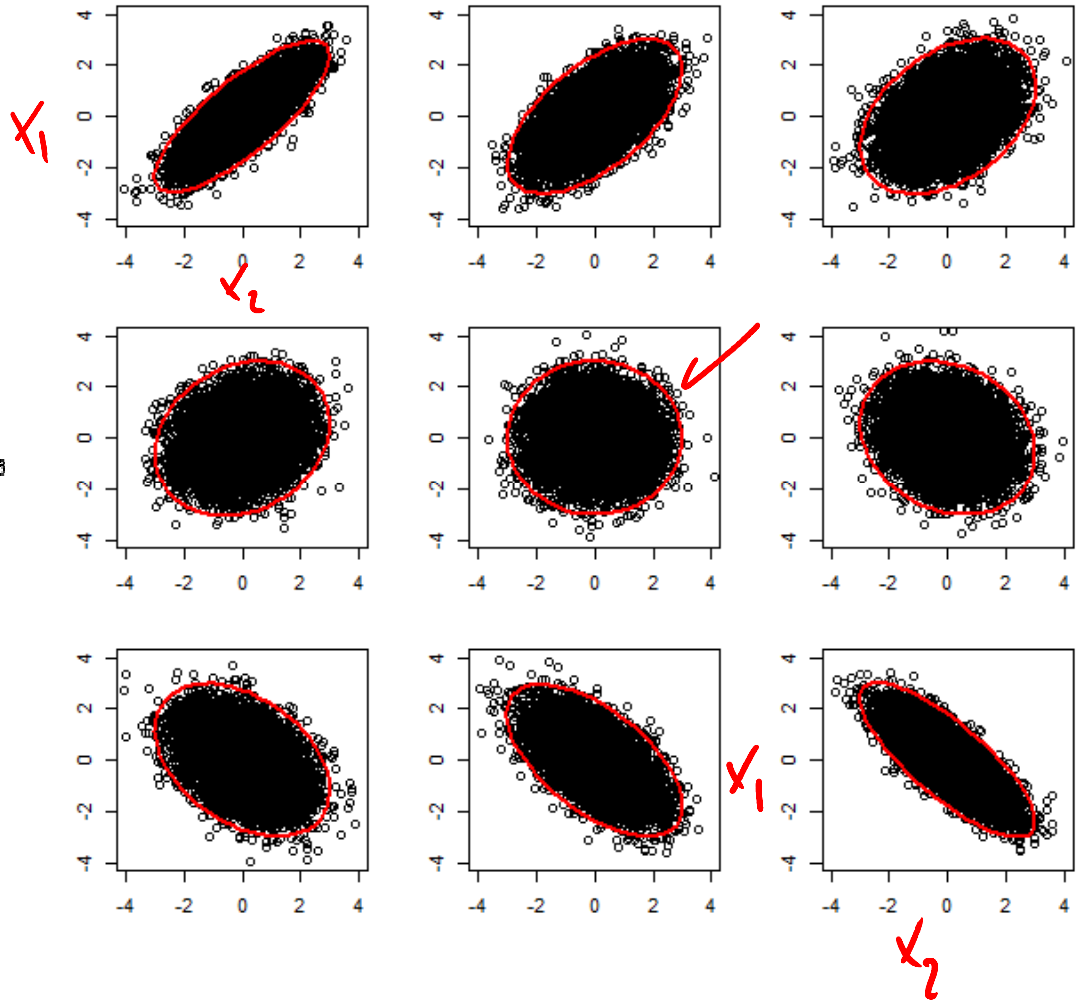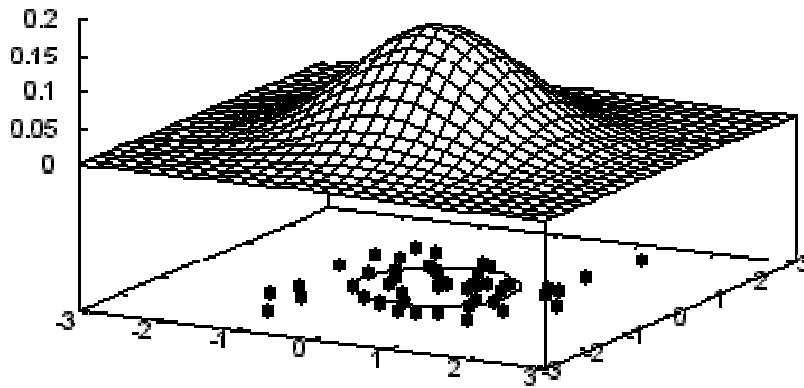$$x \sim \mathcal{N}(\mu, \sigma^2)$$

simulated sample generated

# The bivariate Gaussian distribution

# Multivariate Gaussian distribution

Let $\mathbf{y} \in \mathbb{R}^{n \times 1}$, then pdf of an n-dimensional Gaussian is given by

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

where

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} = \begin{pmatrix} \mathbb{E}(y_1) \\ \vdots \\ \mathbb{E}(y_n) \end{pmatrix}$$

and

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{11} \cdots \sigma_{1n} \\ \cdots \\ \sigma_{n1} \cdots \sigma_{nn} \end{pmatrix} = \mathbb{E}[(\mathbf{y}-\boldsymbol{\mu})(\mathbf{y}-\boldsymbol{\mu})^T]$$

$M_1$  $M_2$

$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$  $\Sigma = \begin{bmatrix} 1 & .5 \\ .5 & 1 \end{bmatrix}$  $\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{bmatrix}$

# Bivariate Gaussian distribution example

*Assume we have two **independent** univariate Gaussian variables*

$$x_1 = \mathcal{N}(\mu_1, \sigma^2) \quad and \quad x_2 = \mathcal{N}(\mu_2, \sigma^2)$$

*Their joint distribution* $p(x_1, x_2)$ *is:*

$$\Sigma = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

$$= \sigma^2 I \qquad |\Sigma| = \sigma^4$$

$$P(x_1, x_2) = P(x_1) P(x_2)$$

$$= (2\pi\sigma^2)^{-1/2} e^{-\frac{1}{2\sigma^2}(x_1 - \mu_1)^2} (2\pi\sigma^2)^{-1/2} e^{-\frac{1}{2\sigma^2}(x_2 - \mu_2^2)}$$

$$= (2\pi\sigma^2)^{-\frac{2}{2}} e^{-\frac{1}{2}\left[(x_1-\mu_1)^T(\sigma^2)^{-1}(x_1-\mu_1) + (x_2-\mu_2)^T(\sigma^2)^{-1}(x_2-\mu_2)\right]}$$

$$= (2\pi\sigma^2)^{-1} e^{-\frac{1}{2}\begin{bmatrix}(x_1-\mu_1) & (x_2-\mu_2)\end{bmatrix}\begin{bmatrix}\sigma^2 & 0 \\ 0 & \sigma^2\end{bmatrix}^{-1}\begin{bmatrix}x_1-\mu_1 \\ x_2-\mu_2\end{bmatrix}}$$

$$\underbrace{}_{\Sigma}$$

# Sampling from a multivariate Gaussian distribution

$$X \sim N(\mu, \sigma^2)$$

$$\|\|\|$$

$$X \sim \mu + \sigma N(0,1)$$

$$\underline{X} \sim N(\underline{\mu}, \underline{\underline{\Sigma}}) \quad , \quad \underline{\underline{\Sigma}} = \underline{\underline{B}} \, \underline{\underline{B}}^T \quad (\text{cholesky})$$

$$\underline{X} \sim \underline{\mu} + \underline{\underline{B}} \, N(0, \underline{\underline{I}})$$

$$\mathbf{Y} \sim N(0, 1)$$

*We have **n=3** data points $y_1 = 1$, $y_2 = 0.5$, $y_3 = 1.5$, which are independent and Gaussian with **unknown** mean $\theta$ and variance 1:*

$$y_i \sim \mathcal{N}(\theta, 1) = \theta + \mathcal{N}(0, 1)$$

*with likelihood $P(y_1 y_2 y_3 | \theta) = P(y_1 | \theta) P(y_1 | \theta) P(y_3 | \theta)$. Consider two guesses of $\theta$, 1 and 2.5. Which has higher likelihood?*

*Finding the $\theta$ that maximizes the likelihood is equivalent to moving the Gaussian until the product of 3 green bars (likelihood) is maximized.*

# The likelihood for linear regression

*Let us assume that each label* $\textbf{\textit{y}}_{\textbf{\textit{i}}}$ *is Gaussian distributed with mean* $\textbf{\textit{x}}_{\textbf{\textit{i}}}^{\textbf{\textit{T}}}\boldsymbol{\theta}$
*and variance* $\boldsymbol{\sigma^2}$*, which in short we write as:*

$$y_i = \mathcal{N}(\, x_i^T\theta \,,\, \sigma^2 \,) \;=\; x_i^T\theta + \mathcal{N}(\, 0,\, \sigma^2 \,)$$

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \sigma) = \prod_{i=1}^{n} p(y_i|\mathbf{x}_i, \boldsymbol{\theta}, \sigma).$$

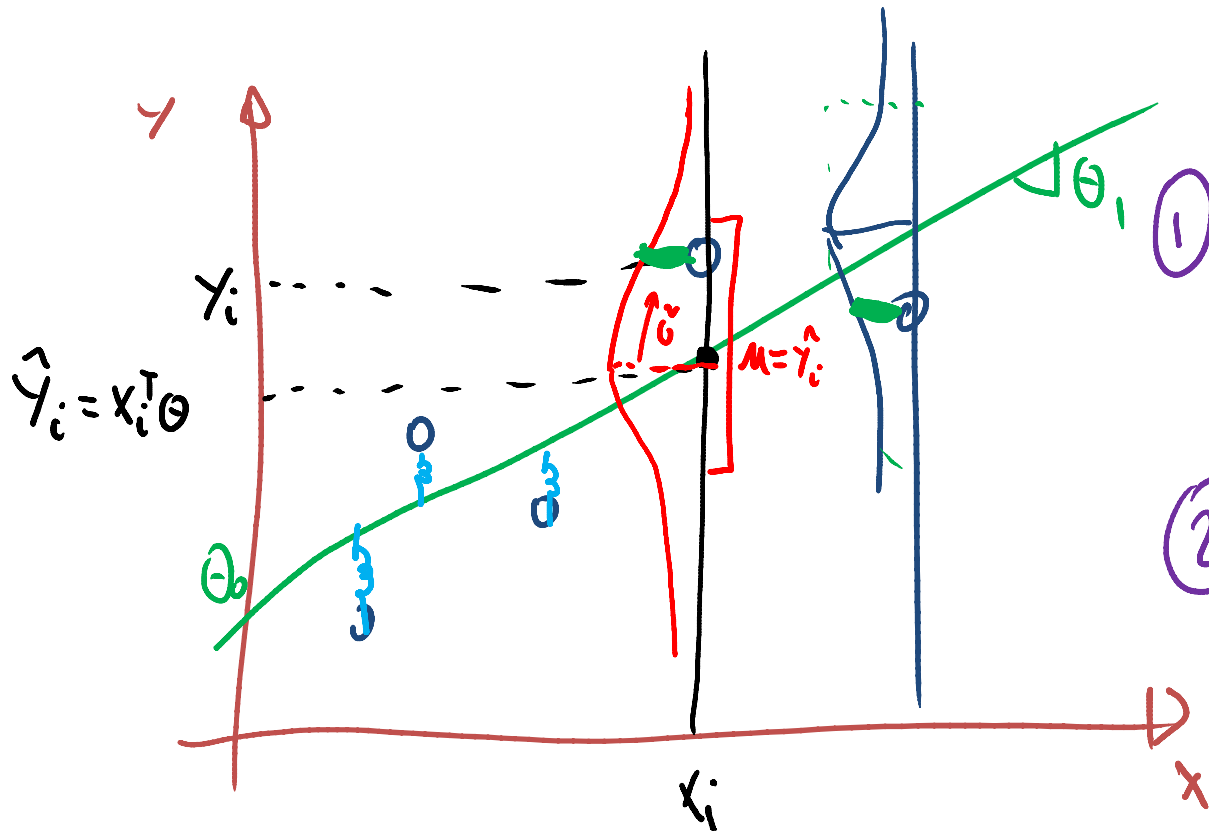$$e^A e^B = e^{A+B}$$

$$= \prod_{i=1}^{n} (2\pi\sigma^2)^{-1/2} e^{-\frac{1}{2\sigma^2}(y_i - \mathbf{x}_i^T\boldsymbol{\theta})^2}$$

$$= (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2}\sum_{i=1}^{n}(y_i - \mathbf{x}_i^T\boldsymbol{\theta})^2}$$

$$= (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})}$$

$$J(\boldsymbol{\theta}) = \sum_{i=1}^{n} (y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2$$

$$p(\mathbf{y} \,|\, \underbrace{\mathbf{X}, \boldsymbol{\theta}, \sigma}_{\text{given}}) = (2\pi\sigma^2)^{-n/2} \, e^{-\frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2}$$

y given

$$= \prod_{i=1}^{n} (2\pi\sigma^2)^{-1/2} \, e^{-\frac{1}{2\sigma^2} (y_i - \mathbf{x}_i^T \boldsymbol{\theta})^2}$$



① Minimize COST

Sum $\left(\frac{b}{i}\right)$

② Maximize PROB.

Product $(\rightarrow)$

# Maximum likelihood

The maximum likelihood estimate (MLE) of $\boldsymbol{\theta}$ is obtained by taking the derivative of the log-likelihood, $\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \sigma)$. The goal is to maximize the likelihood of seeing the training data $\mathbf{y}$ by modifying the parameters $(\boldsymbol{\theta}, \sigma)$.

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \sigma) = \left(2\pi\sigma^2\right)^{-n/2} e^{-\frac{1}{2\sigma^2}(\mathbf{y}-\mathbf{X}\boldsymbol{\theta})^T(\mathbf{y}-\mathbf{X}\boldsymbol{\theta})}$$

$$\ell(\theta) = -\frac{n}{2}\log\left(2\pi\sigma^2\right) - \frac{1}{2\sigma^2}(y-X\theta)^T(y-X\theta)$$

*The ML estimate of $\boldsymbol{\theta}$ is:*

$$\ell(\theta) = -\frac{h}{2} \log\left(2\pi\sigma^2\right) - \frac{1}{2\sigma^2}(y - X\theta)^T(Y - X\theta)$$

$$\frac{\partial \ell(\theta)}{\partial \theta} = 0 - \frac{1}{2\sigma^2}\left[0 - 2X^T y + 2X^T X \theta\right]$$

equating to zero

$$\hat{\theta}_{ML} = \left(X^T X\right)^{-1} X^T y$$

*The ML estimate of σ is:*

$$\frac{\partial}{\partial \sigma} \ell(\sigma)$$

$$\vdots \quad \text{DO THIS !}$$
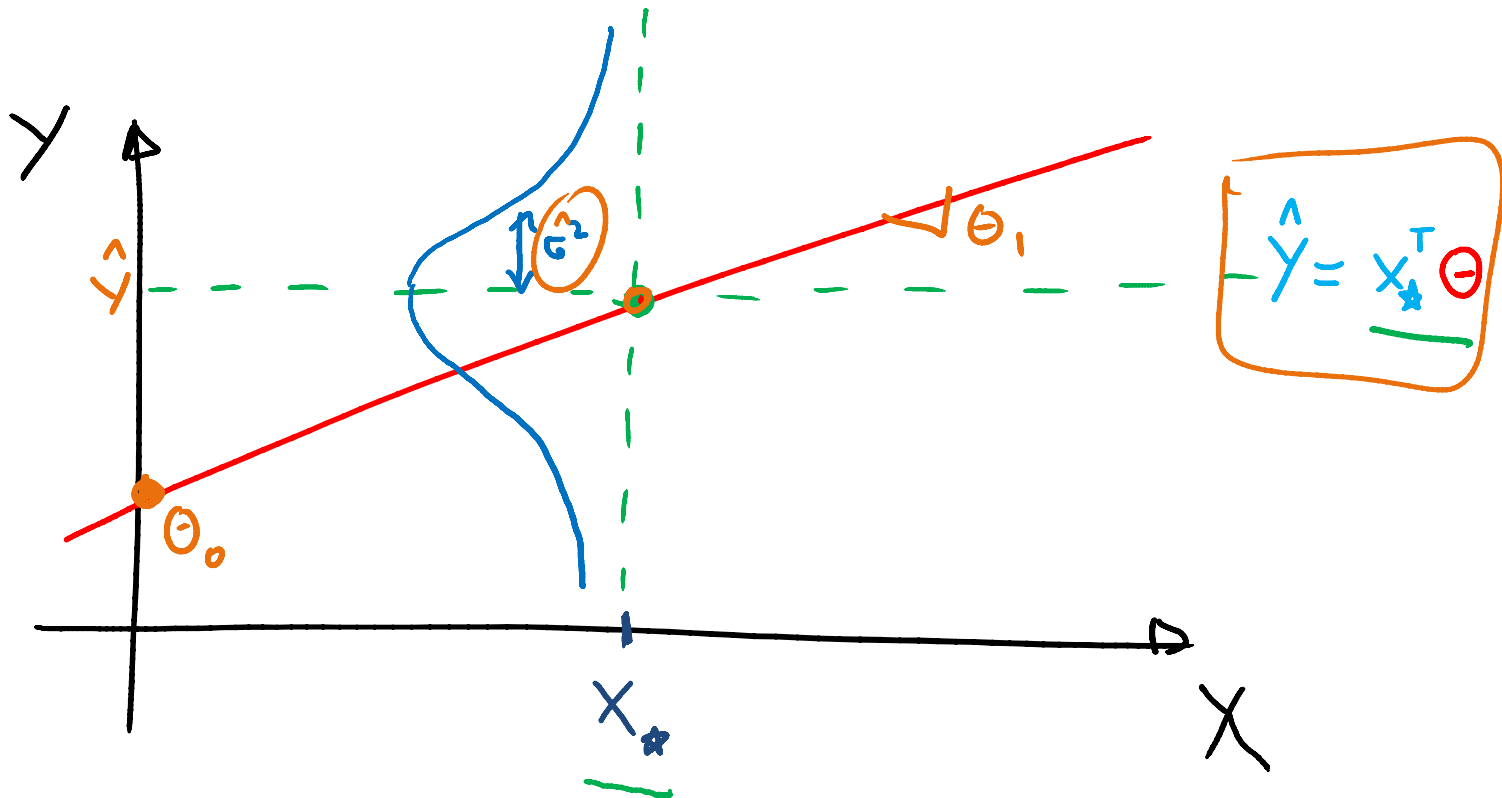
$$\sigma^2 = \frac{1}{n}(Y - X\theta)^T(Y - X\theta) = \frac{1}{n}\sum_{i=1}^{n}(Y_i - X_i\theta)^2$$

# Making predictions

*The ML plugin prediction, given the training data $D=(X, y)$, for a new input $x_*$ and known $\sigma^2$ is given by:*

$$\to \hat{\theta}_{ML} \quad \hat{\sigma}_{ML}$$

$$P(y \mid x_*, D, \hat{\sigma}^2) = \mathcal{N}(y \mid x_*^T \theta_{ML}, \hat{\sigma}^2)$$

# Frequentist learning and maximum likelihood

*Frequentist learning assumes that there exists a true model, say with parameters $\theta_o$.*

*The estimate (learned value) will be denoted $\hat{\theta}$.*

*Given $n$ data, $x_{1:n} = \{x_1, x_2, ..., x_n\}$, we choose the value of $\theta$ that has more probability of generating the data. That is,*

$$\hat{\theta} = \underset{\theta}{arg\ max}\ \ p(\,x_{1:n}\,/\theta\,)$$

# Bernoulli: a model for coins

*A **Bernoulli random variable r.v. X** takes values in {0,1}*

$$p(x/\theta) = \begin{cases} \theta & \text{if} \quad x=1 \\ 1-\theta & \text{if} \quad x=0 \end{cases}$$

$p(x|\theta)$

*Where $\theta \in (0,1)$. We can write this probability more succinctly as follows:*

$$P(x|\theta) = \theta^x (1-\theta)^{1-x} = \begin{cases} \theta & x=1 \\ 1-\theta & x=0 \end{cases}$$

# Entropy

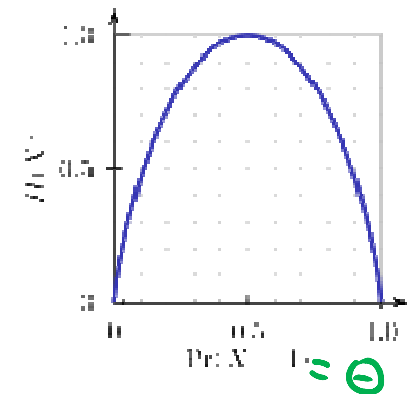*In information theory, entropy **H** is a measure of the uncertainty associated with a random variable. It is defined as:*

$$H(X) = -\sum_{x} p(x|\theta) \log p(x|\theta)$$

*Example:* *For a Bernoulli variable **X**, the entropy is:*

$$H(x) = -\sum_{x=0}^{1} \theta^x (1-\theta)^{1-x} \log\left[\theta^x (1-\theta)^{1-x}\right]$$

$$= -\left[(1-\theta) \log(1-\theta) + \theta \log\theta\right]$$

# MLE - properties

For independent and identically distributed (i.i.d.) data from $p(x|\boldsymbol{\theta}_0)$, the MLE minimizes the **Kullback-Leibler divergence**:

$$\hat{\boldsymbol{\theta}} = \arg\max_{\boldsymbol{\theta}} \prod_{i=1}^{N} p(x_i|\boldsymbol{\theta}) = \arg\max_{\boldsymbol{\theta}} P(X_{1:n}|\theta) \quad \text{given}$$

$$= \arg\max_{\boldsymbol{\theta}} \sum_{i=1}^{N} \log p(x_i|\boldsymbol{\theta}) = \arg\max_{\boldsymbol{\theta}} \log P(X_{1:n}|\theta)$$

$$= \arg\max_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^{N} \log p(x_i|\boldsymbol{\theta}) - \frac{1}{N} \sum_{i=1}^{N} \log p(x_i|\boldsymbol{\theta}_0)$$

$$= \arg\max_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^{N} \log \frac{p(x_i|\boldsymbol{\theta})}{p(x_i|\boldsymbol{\theta}_0)} \qquad X_i \sim P(X_i|\theta_0)$$

$$\xrightarrow{N \to \infty} \arg\min_{\boldsymbol{\theta}} \int \log \frac{p(x|\boldsymbol{\theta}_0)}{p(x|\boldsymbol{\theta})} p(x|\boldsymbol{\theta}_0) dx \simeq \text{KL divergence}$$

LLN.

$$E(x) = \int x \, p(x) \, dx$$

$$N \to \infty$$

$$\hat{E}(x) = \frac{1}{N} \sum_{i=1}^{N} x_i$$

True
when

$$X_i \sim P(x)$$

$$N \to \infty$$

# MLE - properties

$$\arg\min_{\boldsymbol{\theta}} \int \log \frac{p(x|\boldsymbol{\theta}_0)}{p(x|\boldsymbol{\theta})} p(x|\boldsymbol{\theta}_0) dx$$

$$= \arg\min_{\theta} \underbrace{\int P(x|\theta_0) \log P(x|\theta_0) dx}_{\substack{\text{information} \\ \text{world}}} - \underbrace{\int P(x|\theta) \log P(x|\theta) dx}_{\substack{\text{in} \\ \text{model}}}$$

# MLE - properties

Under smoothness and identifiability assumptions, the MLE is **consistent**:

$$\hat{\boldsymbol{\theta}} \xrightarrow{p} \boldsymbol{\theta}_0$$

or equivalently,

$$\mathrm{plim}(\hat{\boldsymbol{\theta}}) = \theta_0$$

or equivalently,

$$\lim_{N \to \infty} P(|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0| > \alpha) \to 0$$

for every $\alpha$.

# MLE - properties

The MLE is **asymptotically normal**. That is, as $N \to \infty$, we have:

$$\left(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\right) \Longrightarrow N(0, I^{-1}) \qquad \text{All of Statistics}$$

where $I$ is the **Fisher Information matrix**.

It is asymptotically optimal or **efficient**. That is, asymptotically, it has the lowest variance among all well behaved estimators. In particular it attains a lower bound on the CLT variance known as the **Cramer-Rao lower bound**.

But what about issues like robustness and computation? Is MLE always the right option?

# Bias and variance

Note that the estimator is a function of the data: $\boxed{\hat{\boldsymbol{\theta}} = g(\mathcal{D})}$ $D = X_{1:n}$
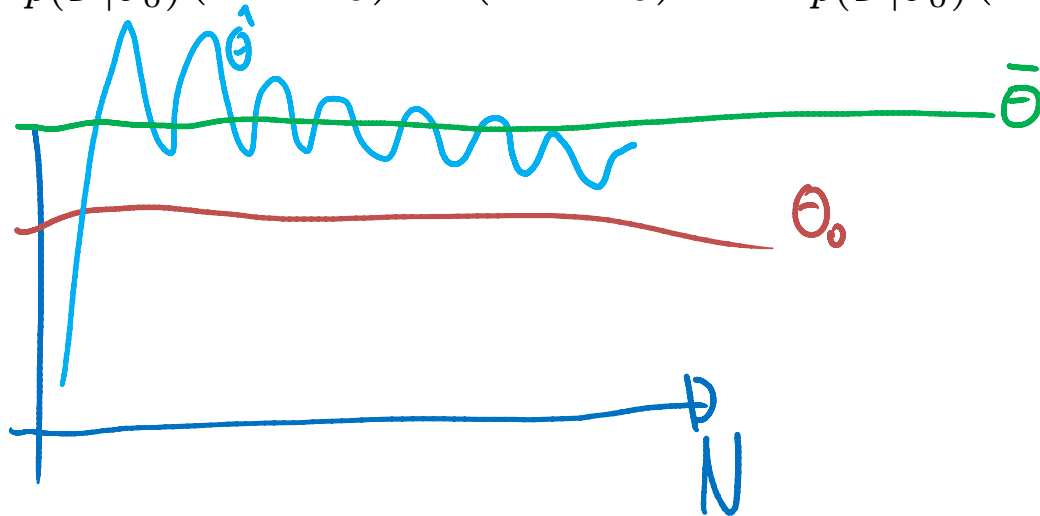
Its **bias** is:

$$bias(\hat{\boldsymbol{\theta}}) = \mathbb{E}_{p(\mathcal{D}|\boldsymbol{\theta}_0)}(\hat{\boldsymbol{\theta}}) - \boldsymbol{\theta}_0 = \bar{\boldsymbol{\theta}} - \boldsymbol{\theta}_0$$

$$\bar{\Theta} = \int \hat{\Theta} \, P(D|\theta_0) \, dD$$

Its **variance** is:

$$\mathbb{V}(\hat{\boldsymbol{\theta}}) = \mathbb{E}_{p(\mathcal{D}|\boldsymbol{\theta}_0)}(\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}})^2$$

Its **mean squared error** is:

$$\text{MSE} = \mathbb{E}_{p(\mathcal{D}|\boldsymbol{\theta}_0)}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^2 = (\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^2 + \mathbb{E}_{p(\mathcal{D}|\boldsymbol{\theta}_0)}(\hat{\boldsymbol{\theta}} - \bar{\boldsymbol{\theta}})^2$$

# Next lecture

In the next lecture, we introduce ridge regression and the Bayesian learning approach for linear predictive models.