# Gradient Descent

- Initialization $W \leftarrow W_0$
- Repeat until convergence $\longrightarrow \|\nabla_w E\| < \epsilon$

$$W_{t+1} \leftarrow W_t - \eta \overbrace{\nabla_w E} \longrightarrow E(D, W) = \sum_{i=1}^{n} \left( W^T x_i - y_i \right)^2$$

$\underset{\text{learning rate}}{\eta}$

$$\nabla_w E = \sum_{i=1}^{n} \nabla_w E_i$$

GD is excellent in accuracy
expensive in computation.

# Stochastic GD

update step: $\quad W_{t+1} \leftarrow W_t - \eta \, \nabla_W E(W, x_i, y_i)$

$i$ is randomly chosen

## fast algorithm

GD     mini-batch GD     SGD

$[n] := \{1, 2, \ldots, n\}$ $\qquad E_i = (W^T x_i - y_i)^2$

all data points $\qquad$ random one point

### MB-GD

update step: $W_{t+1} \leftarrow W_t - \eta \sum_{i \in B} \nabla_W E_i$

$E\left(W, (x_i, y_i)_{i \in B}\right)$

$B \subseteq \{1, \ldots, n\}$

# MLE : Maximum likelihood estimate

$$D = \left\{ (x_i, y_i)_{i \in [n]} \right\}$$

$$W = \Theta$$

$$y_i = w^T x_i + \epsilon_i$$

$$\underset{\theta}{\text{argmax}} \; \underbrace{P(D \mid \theta)}_{\text{likelihood function}} = \Theta_{MLE}$$

## Coin-toss example :

A coin is tossed $n$ times, $y_j$ is the $j^{Th}$ outcome $y_j$ is a Bernoulli RV $\begin{cases} = 1 & \text{w.p. } \theta \\ 0 & \text{w.p. } (1-\theta) \end{cases}$
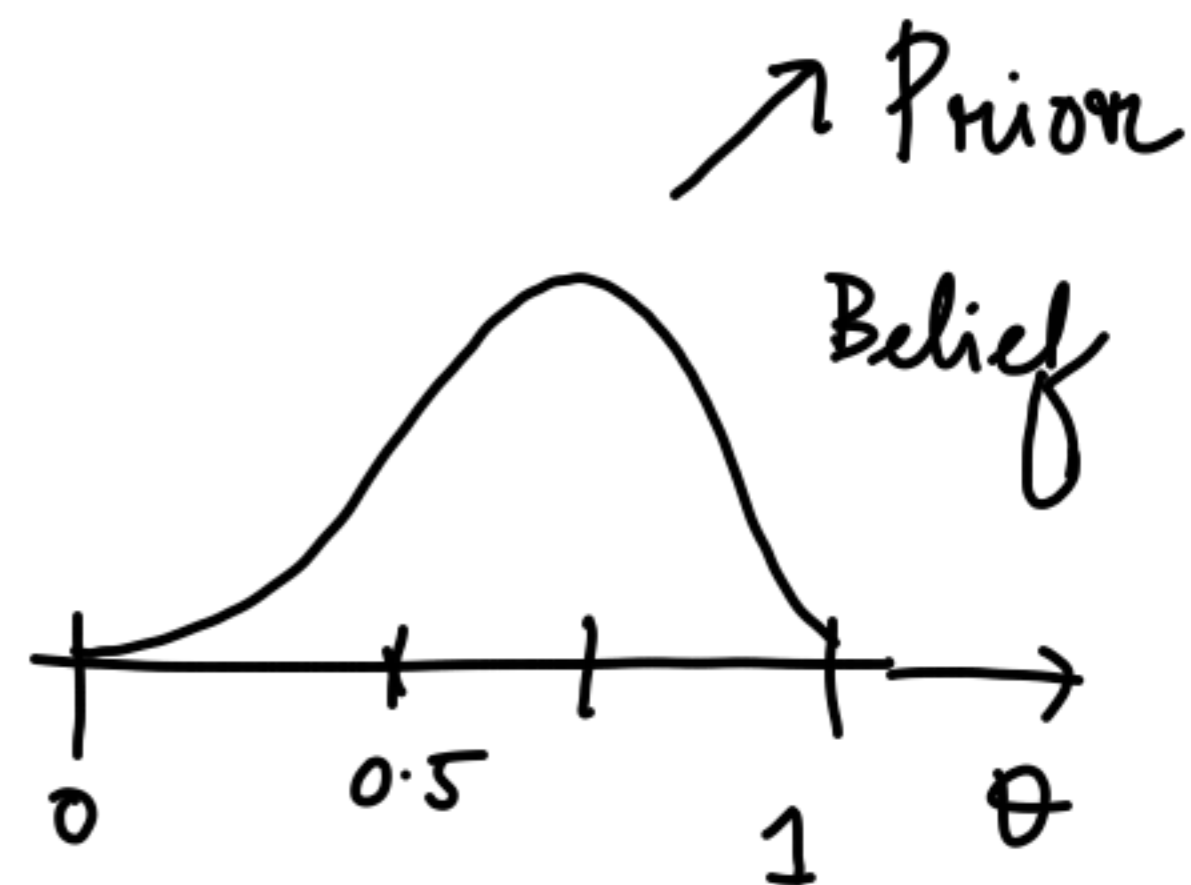
$$P(y_j|\theta) = \theta^{y_j}(1-\theta)^{1-y_j}$$

$$P(y|\theta) = \prod_{i=1}^{n} P(y_i|\theta)$$

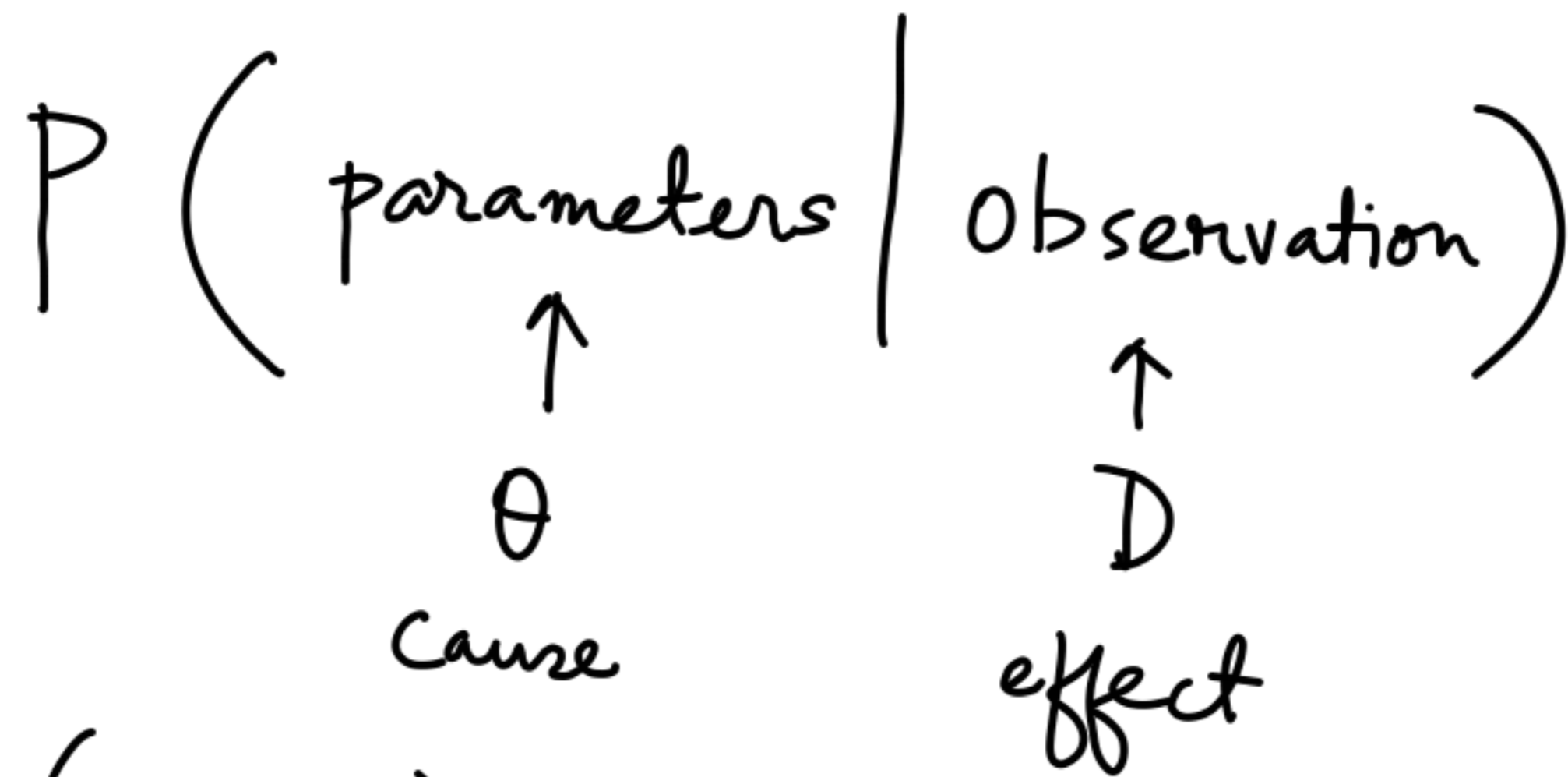$$\underbrace{P(y|\theta)}_{\text{Likelihood}}$$


Prior Belief

$$LL(\theta) = \sum_{i=1}^{n} \log P(y_i|\theta) \implies \theta_{MLE} = \frac{1}{n}\sum_{j=1}^{n} y_j$$

$$P(D|\theta) \qquad \underline{P(\theta)}$$

# Maximum Aposteriori Estimate (MAP)

$$P\left(\underbrace{parameters}_{\substack{\theta \\ cause}} \Big| \underbrace{observation}_{\substack{D \\ effect}}\right)$$

$$P(\theta|D) = \frac{\overset{likelihood}{P(D|\theta)} \overset{prior}{P(\theta)}}{P(D)}$$

Bayesian inference

$P(\theta|D)$ posterior belief

$$\theta_{MAP} \in \underset{\theta}{argmax}\, P(\theta|D) = \underset{\theta}{argmax}\left[P(D|\theta)\, P(\theta)\right]$$

$$\log P(\theta|D) = \log P(D|\theta) + \log P(\theta)$$

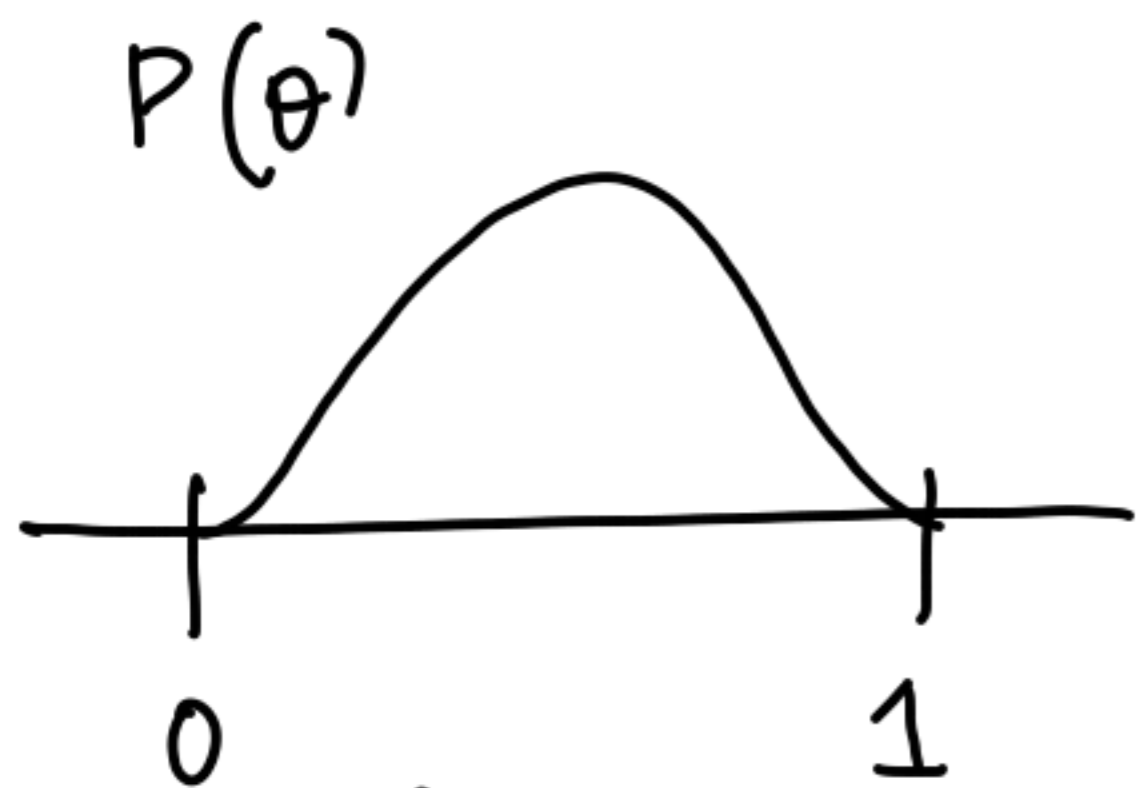$$\theta_{MAP} \in \underset{\theta}{\text{argmax}} \left[ \log P(D|\theta) + \log \underline{P(\theta)} \right]$$

$$\theta_{MAP} = \theta_{MLE} \text{ if } P(\theta) \text{ is constant.}$$
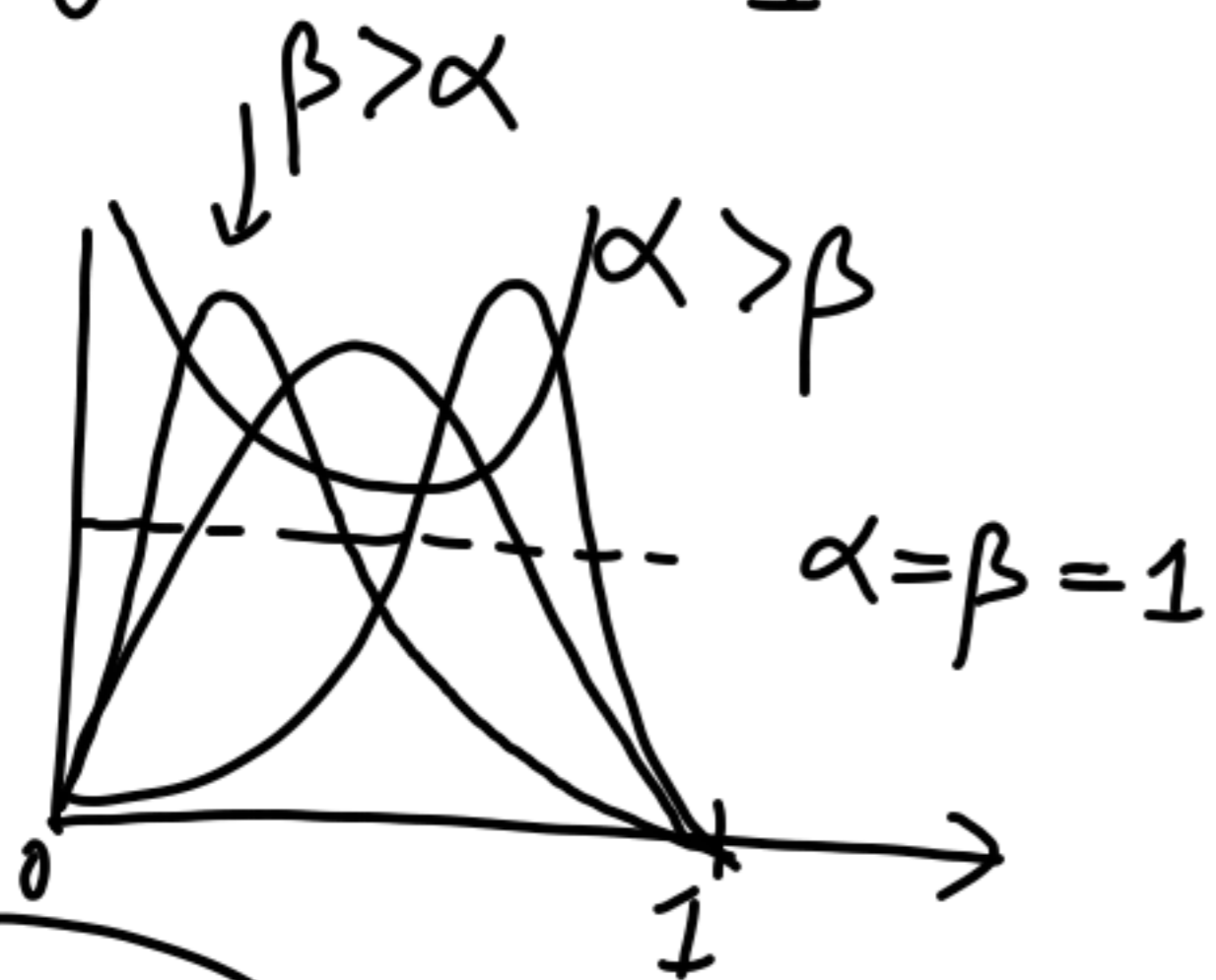
Ex. Likelihood of observing $k$ heads in $n$ tosses

$$P(D|\theta) = \binom{n}{k} \theta^k (1-\theta)^{n-k} \quad \text{Bin}(n,k)$$

$$\theta_{MLE} = \frac{k}{n}$$

$P(\theta)$



Beta $(\alpha, \beta)$

$$P(\theta) = \frac{1}{C} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

- Beta includes a large family of distributions in $[0,1]$

- Beta is <u>conjugate prior</u> of binom dist.

$P(D|\theta) \sim d_1 \quad P(\theta) \sim d_2$

$P(\theta)$ is a CP of $P(D|\theta)$ if $P(\theta|D) \sim d_2$

$P(\theta|D) \propto \underline{P(D|\theta)} \, \underline{P(\theta)}$

$$P(\theta|D) \propto \overbrace{\underbrace{\theta^k (1-\theta)^{n-k}}_{P(D|\theta)}}\ \overbrace{\underbrace{\theta^{\alpha-1} (1-\theta)^{\beta-1}}_{P(\theta)}}$$

$$\propto \theta^{\underbrace{k+\alpha-1}} (1-\theta)^{\underbrace{n-k+\beta-1}} \sim Beta\left(k+\alpha,\ n-k+\beta\right)$$

$$\theta_{MAP} = \underset{\theta}{argmax}\ P(\theta|D) = argmax\ \underbrace{log\ P(\theta|D)}$$

$$= \underset{\theta}{argmax}\ \left[Const. + (k+\alpha-1) log\ \theta + (n-k+\beta-1) log\ (1-\theta)\right]$$

$$\theta_{MAP} = \boxed{\frac{k+\alpha-1}{n+\alpha+\beta-2}}$$

# Conjugate prior examples

1. Bernoulli / Binomial $\longleftrightarrow$ Beta
2. Geometric $\longleftrightarrow$ Beta
3. Categorical $\longleftrightarrow$ Dirichlet
4. · · ·
5. Normal $\longleftrightarrow$ normal $\Big\}$

$$P(\theta | D) \propto P(D | \theta) \, P(\theta)$$

# Conjugate prior for (univariate) Gaussian with known

variance

likelihood $P(D|\theta) \sim N(\mu, \sigma^2)$     (known $\sigma^2$)

prior $P(\theta) \sim N(\mu_0, \sigma_0^2)$

$$\left( \begin{array}{l} D = \{x_1, \ldots, x_n\} \\ = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^{n} (x_i - \mu)^2\right\} \end{array} \right.$$

$$= \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp\left\{-\frac{1}{2\sigma_0^2} (\mu - \mu_0)^2\right\}$$

$$P(\theta|D) \propto P(D|\theta) P(\theta) \propto \exp\left\{-\frac{1}{2\sigma^2} \sum(x_i - \mu)^2 - \frac{1}{2\sigma_0^2} (\mu - \mu_0)^2\right\}$$

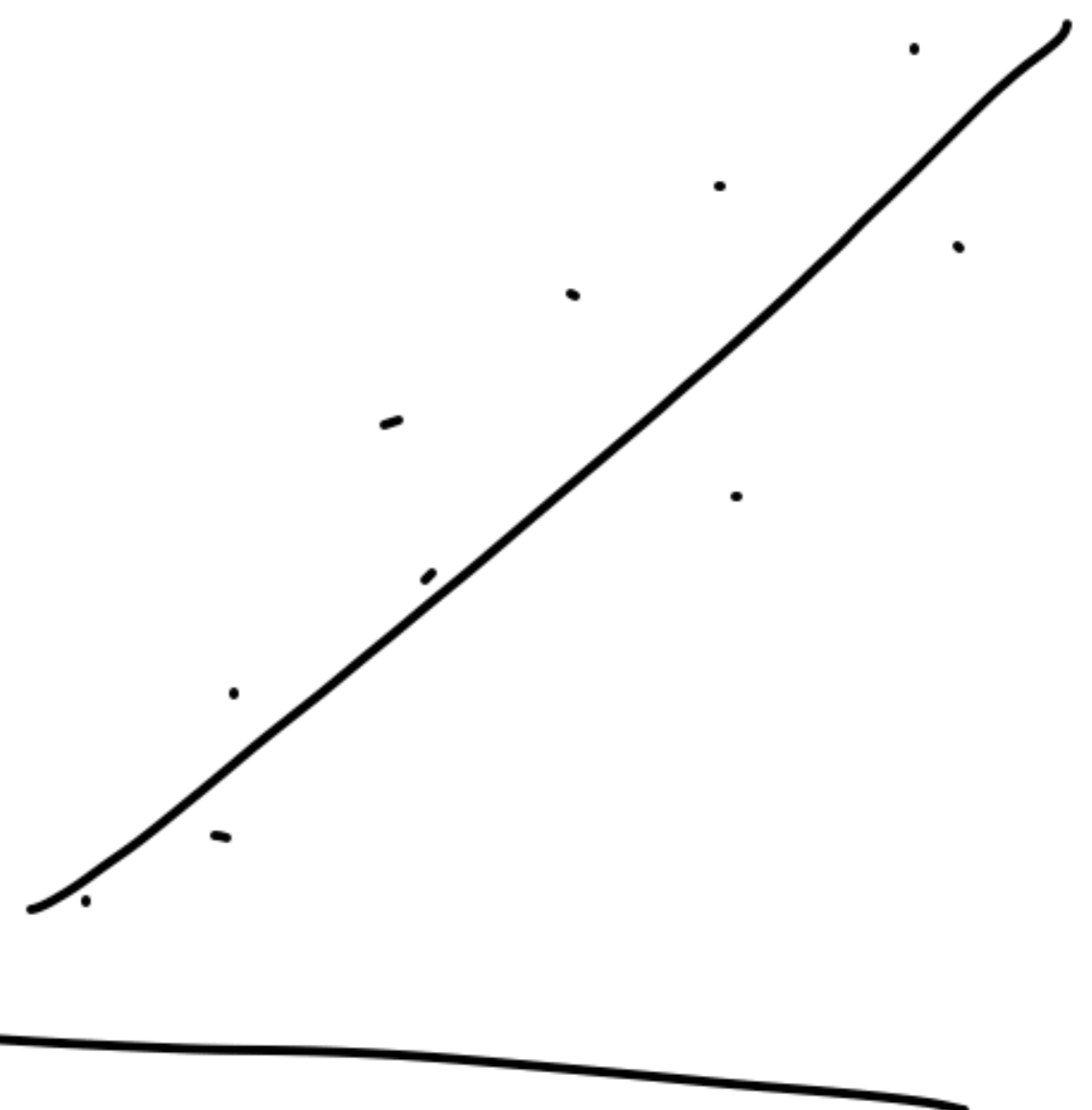Find $\tilde{\sigma}^2$ and $\tilde{\mu}$    $\propto \exp\left\{-\frac{1}{2\tilde{\sigma}^2} \sum(x_i - \tilde{\mu})^2\right\}$

# MAP estimate for linear regression

$$N(w^T x_i, \sigma^2) \leftarrow \underline{y_i} = w^T x_i + \underline{\epsilon_i} \sim N(0, \sigma^2)$$

$$P(D|\theta) \propto \exp\left\{ -\frac{1}{2\sigma^2} \sum (y_i - w^T x_i)^2 \right\}$$

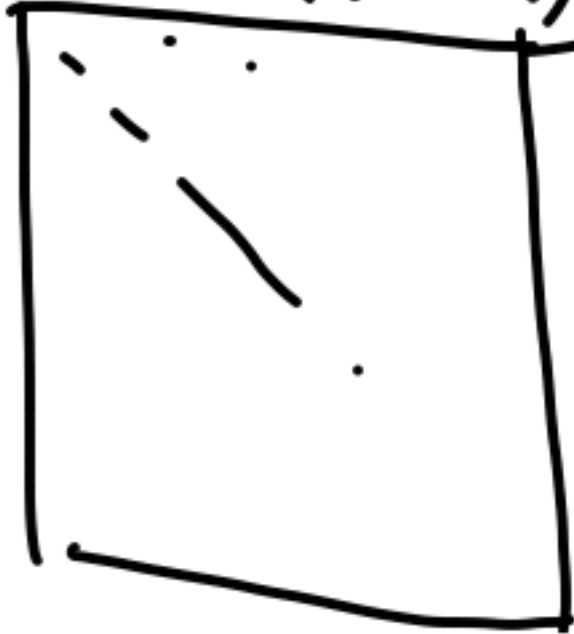$$\theta_{MLE} = \underset{w}{\text{argmax}} \sum_{i=1}^{n} (y_i - w^T x_i)^2$$

$$\frac{P(w)}{P(x_1, x_2, \cdots, x_n)} \sim N\left(0, \frac{1}{\lambda} I\right)$$

$$\lambda > 0$$

$$P(x) \sim N(\mu, \underset{\downarrow}{\Sigma})$$

$$x \in \mathbb{R}^d$$

$$d \times d$$

$$\sigma_{ij} = \text{cov}(x_i, x_j)$$

$$E[(x_i - Ex_i)(x_j - Ex_j)]$$

$$P(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right\}$$

multivariate normal distribution.

$$P(w) = \frac{1}{(2\pi)^{d/2} \left(\frac{1}{\lambda}\right)^{d/2}} \exp\left\{-\frac{\lambda}{2} w^T w\right\}$$

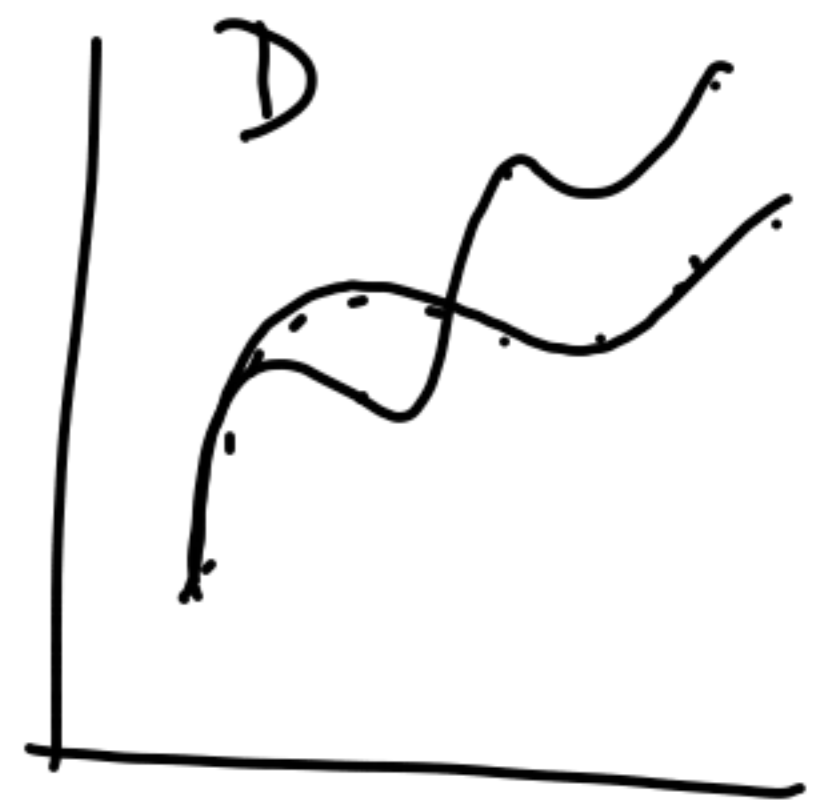$$\sim N\left(0, \frac{1}{\lambda}I\right)$$

$$\propto \exp\left\{-\frac{\lambda}{2}\|w\|^2\right\}$$

$$P(w|D) \propto P(D|w) \, P(w)$$

$\lambda = \text{hyperparameter}$

$$\arg\max_{w} \left[ \log P(D|w) + \log P(w) \right]$$

$$\arg\min_{w} \left\{ \underbrace{\frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - w^T x_i)^2}_{\hat{y}} + \frac{\lambda}{2} \|w\|^2 \right\}$$

$$\arg\min_{w} \left\{ \frac{1}{2\sigma^2} \|Xw - y\|^2 + \underbrace{\frac{\lambda}{2} \|w\|^2}_{\text{Regularizer}} \right\}$$

$w_0 + w_1 x + w_2 x^2 + \dots +$