Proba supo learning & MLE $\Longleftrightarrow$ ERM

ERM: $R(\theta)$

RRM: $R(\theta) = PredLoss(\theta) + Regul(\theta)$

$$\lambda \|\theta\|^2$$

Proba. inferp as Bayesian learning

---

MLE vs. (map) Bayesian estim. $\overbrace{Lik(\theta)}$

$: Z, \theta$          $\hat{\theta}_{MLE} = \underset{\theta}{argmax} \; P(z|\theta)$

$\underset{data}{\uparrow} \; \underset{param}{\uparrow}$



$\hat{\theta}_{MLE}$

Bayesian estim.

prior distrib   $P(\theta|\lambda)$ $\longleftarrow$ fixed hyperparam.

$\leftarrow$ prior knowledge

$\theta \hookrightarrow \underline{\underline{RV}} !$

$\lambda \rightarrow \theta \rightarrow Z$

Joint distrib:  $P(z, \theta|\lambda) = P(z|\theta, \cancel{\lambda}) \, P(\theta|\lambda)$

Posterior distrib. of $\theta$

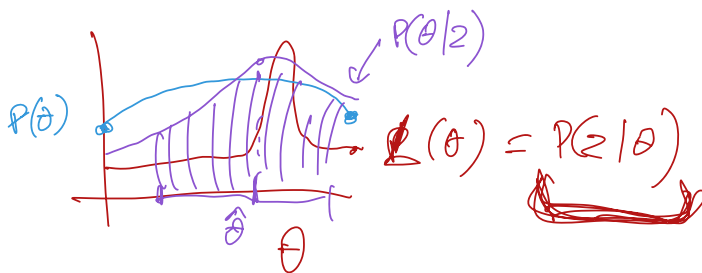$$P(\theta|z, \lambda) = \frac{P(z|\theta) \, P(\theta|\lambda)}{P(z|\lambda)}$$

$P(\lambda)$

$$P(\theta \mid z, \lambda) = \frac{1}{P(z \mid \lambda)} P(z \mid \theta) \, P(\theta \mid \lambda)$$

$$\propto \qquad P(z \mid \theta) \, P(\theta \mid \lambda) \qquad \text{``Unnorm. Posterior of } \theta \text{''}$$

## MAP = Maximum a Posteriori

$$\hat{\theta}_{MAP} = \underset{\theta}{\text{argmax}} \; P(z \mid \theta) \, P(\theta \mid \lambda)$$

if $P(\theta \mid \lambda) = $ const, $\implies \hat{\theta}_{MAP} = \hat{\theta}_{MLE}$



$P(\theta \mid z)$

$P(\theta)$

$\mathcal{L}(\theta) = P(z \mid \theta)$

$\hat{\theta}$

$\theta$

Prob. posterior for $\theta$ : More stuff
- = MAP ]
- = Mean
- = 95% CI

# L2-Reg. LinReg as Generative Model

$X, \sigma^2$: fixed $\qquad \vec{\theta} \; \vec{y} : R.V. \qquad \vec{X}_i = \text{fixed}$

① For $j = 1..D$: $\quad P(\theta_j | \lambda) = N\left(0, \frac{1}{\lambda}\right) \quad \Longleftrightarrow \quad \theta_j \sim N\left(0, \frac{1}{\lambda}\right)$

② for $i = 1.. N$: $\quad P(y_i | x_i, \theta, \sigma^2) = N(\theta^T x_i, \sigma^2)$

## Joint Prob. Distrib

$$P(\theta, Y | X, \sigma^2, \lambda)$$
$$= P(Y | X, \theta, \sigma^2, \cancel{\lambda}) \; P(\theta | \lambda, \cancel{X}, \cancel{\sigma^2})$$

## Dir. Graphical Model



◎ Observed/fixed/conditioned-on
○ Latent

$$\hat{\theta}_{MAP} = \underset{\theta}{\text{argmax}} \ P(\theta \mid Y, X, \sigma^2, \lambda)$$

$$P(Y \mid X, \sigma^2, \theta) \quad P(\theta \mid \lambda)$$

$$\log P(Y \mid X, \sigma^2, \theta) + \log P(\theta \mid \lambda) \qquad N\left(0, \sigma^2 = \frac{1}{\lambda}\right)$$

Prior

$$\log P(\theta \mid \lambda) = \sum_{j=1}^{D} \log\left[ \frac{1}{\sqrt{2\pi \frac{1}{\lambda}}} \exp\left[ \frac{-1}{2\left(\frac{1}{\lambda}\right)} \theta_j^2 \right] \right]$$

$$= \sum_j \left[ \underbrace{\left(-\log \sqrt{2\pi}\right) + \log \sqrt{\lambda}}_{\text{constant } C} - \frac{\lambda}{2} \theta_j^2 \right]$$

$$= C + \frac{\lambda}{2} \sum_j \theta_j^2$$

<u>Like</u> $\quad P(Y/X, \theta, \sigma^2) = \prod_i \underbrace{P(y_i/x_i, \theta, \delta^2)}_{N(\theta^T x_i, \sigma^2)}$

$\log(Y/X, \theta, \sigma^2) =$

$\sum_i \log\left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-1}{2\sigma^2}\left(y_i - \theta^T x_i\right)^2\right)\right]$

$= \sum_i \underbrace{\log\left[\frac{1}{\sqrt{2\pi\sigma^2}}\right]}_{\text{const wrt } \theta} - \frac{1}{2\sigma^2}\left(y_i - \theta^T x_i\right)^2$

$= -\frac{1}{2\sigma^2} \sum_i \left(y_i - \theta^T x_i\right)^2 \qquad + \text{Const}$

$J(\theta) =$

$\log P(Y/X, \theta, \sigma^2) + \log P(\theta/\lambda)$

$= -\frac{1}{2\sigma^2} \sum_i \left(y_i - \theta^T x_i\right)^2 - \frac{\lambda}{2} \sum_j \theta_j^2 \qquad + \text{const}$

$\hat{\theta}_{MAP} = \underset{\theta}{\text{argmax}} \; J(\theta)$
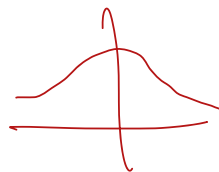
# Why normal?

- Common

- convenient — L2 norm is CVX

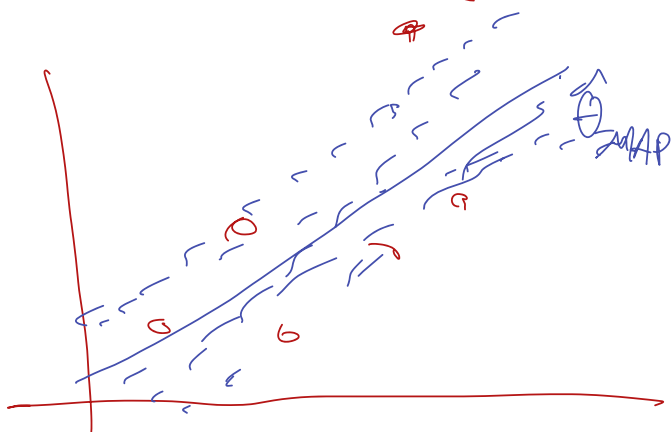# Why param. prob. models?
↳ Makes assumps clear

or: be assump-free ??

---

$$y_i \sim N\left(\sum_j \alpha_j K(x_j, x_i), \, \delta^2\right)$$

$P(\theta \mid Y, X \ldots)$

$\hat{\theta}_{MAP}$