# Tutorial For A3

1)    ridge regression

$$w_{MAP} = \underset{w}{\arg\min} \; \frac{1}{2} \| Xw - Y \|^2 + \frac{\lambda}{2} \|w\|^2$$

↑ minimize objective function     ← don't want $w$ too large
$\lambda$ is large, more important
to minimize thes

think back to MLE:

$$w_{MLE} = \underset{w}{\arg\min} \| Xw - Y \|^2$$

deterministic    gaussian

$$y = x^{\top} w + n, \quad n \sim N(0, \sigma^2)$$

↑ deterministic    noise term probability distribution

data given hypothesis

MLE: $p(D/h)$

MAP: $p(h|D)$

given data, want to find weights

$p(w|y,x)$

(w that would maximize)

$$\propto p(Y|w,x) \, p(w)$$

ex: $n \sim N(0, \sigma^2)$
$p(y|w,x) = N(w^{\top}x, \sigma^2)$

$x_2$

$w_1 x_1 + w_2 x_2$

$x_1$

corresponds to $p(y|w,x)$
corresponds to $p(w)$

$$p(w|y,x) \propto \cdots$$

$$\log(p(w|y,x)) = \log(p(Y|x,w)) + \log(p(w)) + c$$

$$-\log(p(w|y,x)) = -\log(p(Y|x,w)) - \log(p(w)) + c$$

$$-\log(p(w)) \propto \frac{\lambda}{2} \|w\|^2$$

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

to get rid of $\propto$, have to
normalize to sum up $\Sigma p = 1$

$$w \sim N(0, \tfrac{1}{\lambda}I)$$

$$w_i \sim N(0, \tfrac{1}{\lambda})$$

$$\|w\|^2 = w^{\top} I w \quad (w \text{ is vector})$$

$$p(Y|x,w) \propto -\frac{1}{2}\|Xw - Y\|^2$$

$$= k \exp\left(\frac{-1}{2}(Xw-y)^T (xw-y)\right)$$

$$y \sim N(Xw, \mathbb{I}) \qquad \text{vector}$$
$$y_i \sim N(x_i^+w, 1) \qquad \text{elements of the vector}$$

$$y_i \sim N(w^+x_i, 1)$$

$$w_i \sim N(0, \lambda_i) \qquad \Rightarrow \qquad \underset{w}{\text{argmin}} \; \frac{1}{2}\|Xw - Y\|^2 + \frac{1}{2}w^T \Lambda w$$
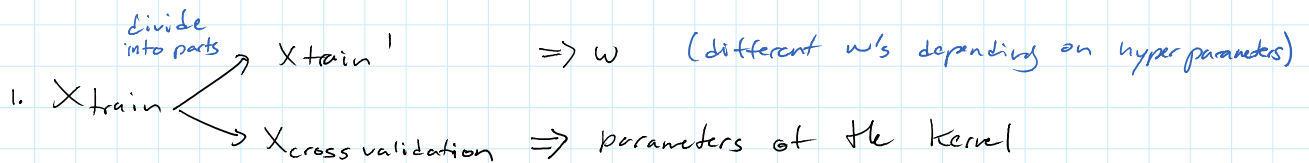
covariance

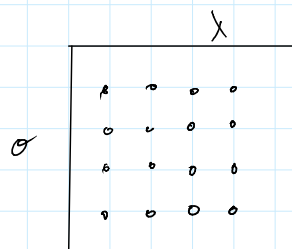$$p(\bar{w}) = k \exp\left(\sum_{i=1}^{d} w_i^2 \left(-\frac{1}{2}\lambda_i\right)\right)$$

$$w^+ \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \; \lambda_d \end{pmatrix} w$$

$$\begin{pmatrix} \lambda_1 w_1 \\ \lambda_2 w_2 \\ \vdots \\ \lambda_d w_d \end{pmatrix}$$

---

3)

1. X train — divide into parts → X train' ⟹ w (different w's depending on hyper parameters)

→ X cross validation ⟹ parameters of the kernel

ii. X test



"hyper parameters"

➔ iterate over $\lambda_i \; \sigma_j$
to find the best combination

X

$\sigma$

**4)**

$$\arg\min_{w} \frac{1}{2} \|Xw - Y\|^2 + \boxed{\lambda \|w\|_1} = w_{MAP}$$

↑ change to L1 norm
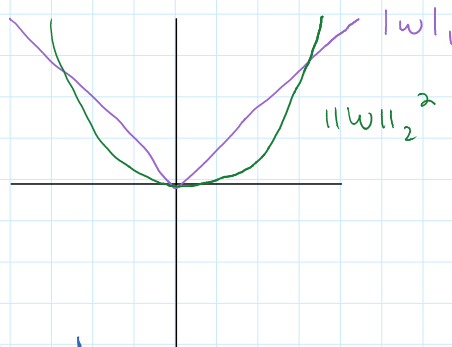
$$p(Y|xw) \quad p(w)$$

→ laplacian (heavy tail)

$$p(w_i) \propto \exp(-\lambda |w_i|)$$

"this is the laplace distribution"

laplace: $f(z|\mu, b) = \frac{1}{2b} \exp\left(\frac{-|z-\mu|}{b}\right)$



$|w|_1$

$\|w\|_2^2$

$$w_i \sim laplace\left(0, \frac{1}{\lambda}\right)$$

$$\int_{-\infty}^{\infty} p(w_i) dw_i = 1$$

upper bound L1 norm

$$\arg\min_{w, v} \frac{1}{2} \|Xw - Y\|^2 + \lambda \sum v_i$$

$$|w_i| < v_i \qquad \text{each feature,}$$
$$w_i < v_i \qquad \text{2 constraints.}$$
$$-w_i < v_i$$

---

**5)**

the algorithm is in the book.

Algorithm 13.1 : Coordinate Descent for lasso

**Algorithm 13.1:** Coordinate descent for lasso (aka shooting algorithm)

1  Initialize $\mathbf{w} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{y}$;
2  **repeat**
3      **for** $j = 1, \ldots, D$ **do**
4          $a_j = 2\sum_{i=1}^{n} x_{ij}^2$;
5          $c_j = 2\sum_{i=1}^{n} x_{ij}(y_i - \mathbf{w}^T\mathbf{x}_i + w_j x_{ij})$;
6          $w_j = \text{soft}(\frac{c_j}{a_j}, \frac{\lambda}{a_j})$;
7  **until** *converged*;