

## Lecture 18, part I

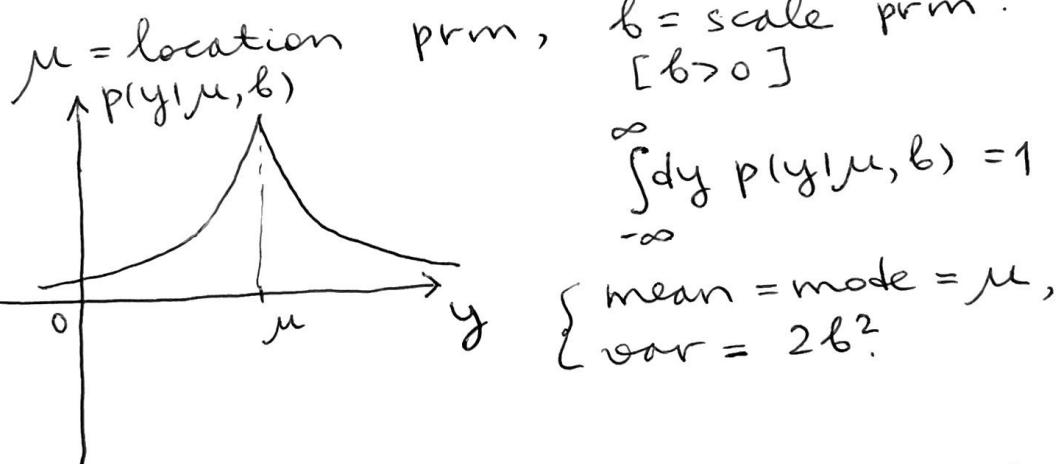
### Robust linear regression

Fitting gaussian models is sensitive to outliers  $\Rightarrow$  one can use heavy-tail distributions to achieve robustness to outliers.

#### ① Laplace likelihood

Laplace distribution:

$$p(y|\mu, b) = \frac{1}{2b} e^{-\frac{|y-\mu|}{b}}$$



Needs techniques like projected grad for robust optimization.

#### ② Student's t-likelihood

Student's t-distribution:

$$p(y|\mu, \sigma, J) = \frac{\Gamma(\frac{J+1}{2})}{\sqrt{\pi J} \Gamma(\frac{J}{2})} \left(1 + \frac{1}{J} \left(\frac{y-\mu}{\sigma}\right)^2\right)^{-\frac{J+1}{2}}$$

$\mu$  = mean,  $\sigma > 0$  = scale prm,  
prm  $J > 0$  = DoF

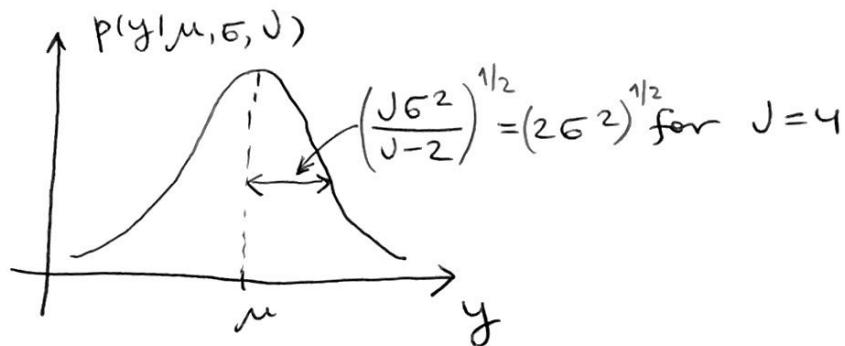
~~mean~~

$$\left\{ \begin{array}{l} \text{mean} = \mu, J > 1 \\ \text{mode} = \mu \\ \text{var} = \frac{J\sigma^2}{J-2}, J > 2 \end{array} \right. \Rightarrow J=4 \text{ is commonly used}$$

Note that for  $t^2 \gg J$ ,

$$p(t) \sim (t^2)^{-\frac{J+1}{2}} = t^{-(J+1)}$$

Thus, the tail decays much slower than the Gaussian tail.

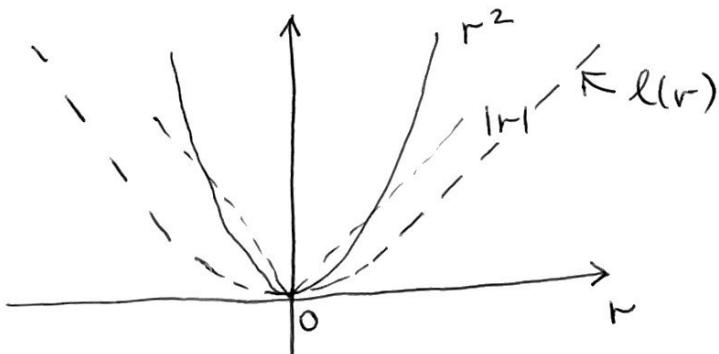


The model can be optimized using SGD or other standard techniques.

### ③ Huber loss

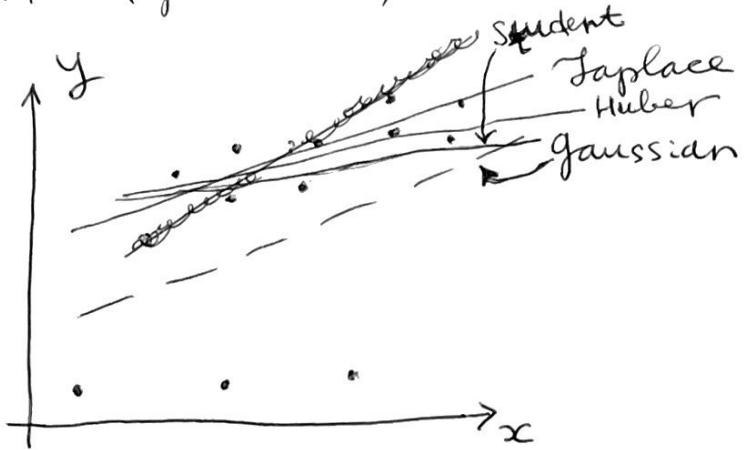
$$l(r) = \begin{cases} r^2/2, & |r| \leq \delta \\ \delta|r| - \delta^2/2, & |r| > \delta \end{cases}$$

fitting weight  
 $r = w - d$   
 constraint,  
 typically = 0,  
 or  
 $r = \text{obs.} - \text{model}$   
 $y - \hat{y}$



Typically,  $\delta = \Theta(1)$ , e.g. 1.5

Huber loss can be used instead of  $|...|^2$  (gaussian) loss or  $|...|$  (laplace) loss.



### Bayesian linear regression

Likelihood:  $p(D|\vec{\omega}, \sigma^2) = \prod_{n=1}^N p(y_n | \vec{\omega} \cdot \vec{x}, \sigma^2) = N(\vec{y} | X\vec{\omega}, \sigma^2 \mathbb{I}_N)$ . centered data

Prior:  $p(\vec{\omega}) = N(\vec{\omega} | \vec{0}, \tau^2 \mathbb{I}_D)$ .

Then the posterior is given by

$p(\vec{\omega} | D, \sigma^2) = N(\vec{\omega} | \vec{\omega}', \Sigma')$ , where

$$\begin{cases} \vec{\omega}' = \frac{1}{\sigma^2} \Sigma' X^\top \vec{y}, \\ \Sigma' = \left( \frac{1}{\tau^2} \mathbb{I}_D + \frac{1}{\sigma^2} X^\top X \right)^{-1} = \sigma^2 \left( \frac{\lambda \mathbb{I}_D + X^\top X}{\tau^2} \right)^{-1}. \end{cases}$$

Then  $\vec{\omega}' = (\lambda \mathbb{I}_D + X^\top X)^{-1} X^\top \vec{y}$ , same as ridge regression.

Finally, we can compute the predictive distribution:

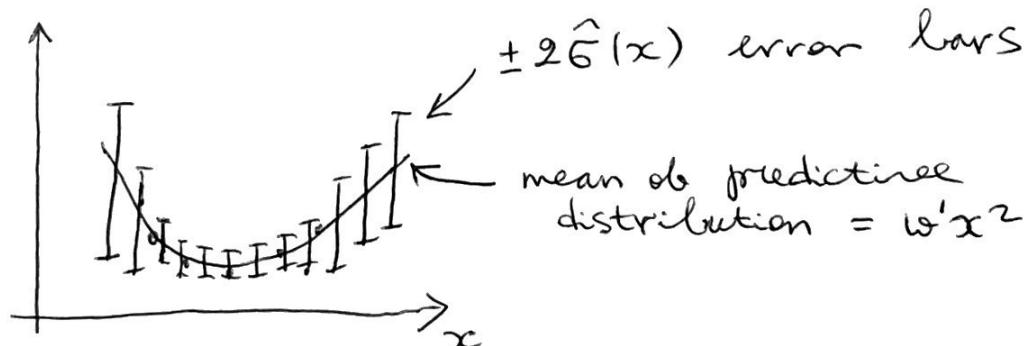
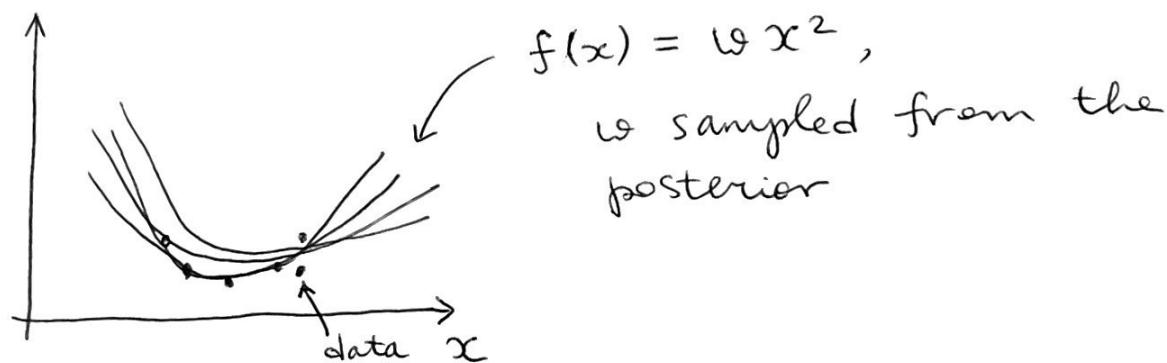
$$p(y|\vec{x}, \mathcal{D}, \sigma^2) = \int d\vec{\omega} N(y|\vec{\omega} \cdot \vec{x}, \sigma^2) \times$$

$$\times N(\vec{\omega} | \vec{\omega}', \Sigma') =$$

$$= N(y | \vec{\omega}' \cdot \vec{x}, \hat{\sigma}^2(\vec{x})) , \text{ where}$$

$$\hat{\sigma}^2(\vec{x}) = \sigma^2 + \vec{x}^\top \Sigma' \vec{x}$$

$\sigma^2$   
 —————  
 ↑  
 var.  
 in observation  
 noise              ↑ var. in posterior  
 N training examples



Now, we switch to features:

$$\tilde{X}_{N \times D} \Rightarrow \Phi_{N \times M}$$

data                  design  
matrix                  matrix

and switch to Bishop's notation:

posterior  $\Rightarrow p(\tilde{\omega} | \mathcal{D}, \sigma^2)$  becomes

$$p(\tilde{\omega} | \mathcal{D}, \lambda, \beta) = \mathcal{N}(\tilde{\omega} | \tilde{m}_N, S_N), \text{ where}$$

$$\begin{cases} \tilde{\omega} \\ \lambda^{-2} \\ \beta^{-2} \end{cases}$$

$$\mathcal{D} = \left\{ \begin{array}{l} \tilde{t}, X \\ \uparrow \text{instead of } \tilde{y} \end{array} \right\} \quad \begin{cases} \tilde{m}_N = \beta S_N \Phi^T \tilde{t}, \\ S_N = (\lambda \mathbb{I} + \beta \Phi^T \Phi)^{-1} \end{cases}$$

Thus, we have relabeled:

$$\begin{cases} \tilde{\omega}' \Rightarrow \tilde{m}_N, \\ \Sigma' \Rightarrow S_N. \end{cases}$$

Finally, the predictive distribution is given by

$$p(t | \tilde{x}, \mathcal{D}, \lambda, \beta) = \mathcal{N}(t | \tilde{m}_N \cdot \underbrace{\tilde{\Phi}(\tilde{x})}_{\substack{\text{M-dim} \\ \text{feature vector}}, \sigma_N^2(\tilde{x})),$$

where

$$\sigma_N^2(\tilde{x}) = \beta^{-1} + \tilde{\Phi}^T(\tilde{x}) S_N \tilde{\Phi}(\tilde{x}).$$

Note that as  $N \rightarrow \infty$ , we expect the posterior prob. to have vanishing variance:  $S_N \rightarrow 0$  as  $N \rightarrow \infty$ .

Indeed,

$$(S_N^{-1})_{ij} \sim \beta (\phi^T \phi)_{ij} =$$

$$= \beta \sum_{n=1}^N \psi_i(\tilde{x}_n) \psi_j(\tilde{x}_n).$$

This  $\uparrow$  as  $N \uparrow \Rightarrow S_N \downarrow$ .

Moreover, if  $\tilde{\psi}(\tilde{x})$  is localized,

$\sigma_N^2(\tilde{x}) \rightarrow \beta^{-1}$  even for small  $N$ ,  
 if  $\tilde{x}$  is far away from the basis  
 function centers. Thus, the model  
 may become overconfident when  
 extrapolating.