

- Review SRM principle

- Bayesian score / BIC

Model Selection

We know a collection of function classes

$$\mathcal{F}_i = \{f(\cdot; \theta) : \theta \in \mathbb{R}^{d_i}, i=1, \dots, K\}$$

Training samples $\mathcal{D}_n = \{\underline{x}_t, y_t\}_{t=1}^n$

Two defn of RISK:

① Expect Risk given parameter $\underline{\theta}$

$$R(\underline{\theta}) = \mathbb{E} [\text{loss}_{0-1}(y; f(x; \underline{\theta}))]$$

$(x,y) \sim P$

joint dist. = $\int \text{Loss}_{0-1}(y; f(x; \underline{\theta})) dP(x, y)$

$X \times \{-1, +1\}$

If we are working in a particular f° class \mathcal{F}_i

$\hat{\theta}_i = \underset{\underline{\theta}: f(\cdot; \underline{\theta}) \in \mathcal{F}_i}{\operatorname{argmin}} J(\underline{\theta}) \Rightarrow$ choose the opt $\hat{\theta}_i$ in every function class \mathcal{F}_i

$$J(\underline{\theta}) = \frac{1}{n} \sum_{t=1}^n \text{Loss}(y_t; f(x_t; \underline{\theta})) + \lambda \|\underline{\theta}\|^2$$



Eg. SVM \rightsquigarrow hinge loss

Logistic Regression \rightsquigarrow logistic loss

Best expected risk within \mathcal{F}_i : $R(\hat{\theta}_i) = R(\hat{f}_i)$, where $\hat{f}_i = \text{fc}(\cdot; \hat{\theta}_i)$

Rmk: • exp. risk can not be computed since $P(x, y)$ is Unknown.

- Approximate the exp. risk with the empirical risk.

(2)

Emp. Risk $\rightarrow \theta$

$$R_n(\theta) = \frac{1}{n} \sum_{t=1}^n \text{Loss}_{\theta-1}(y_t; f(x_t; \theta))$$

$$R_n(\hat{f}_i) = R_n(\hat{\theta}_i) \quad \hat{f}_i \in \mathcal{F}_i$$

↑
Best empirical risk

↓
can be calculated

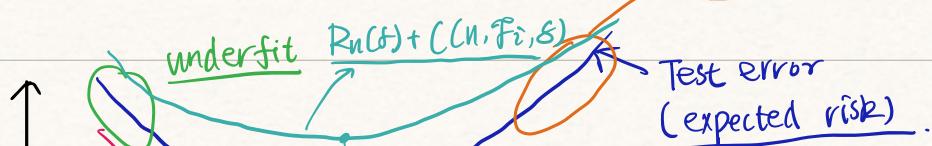
Thm: With probability at least $1-\delta$.

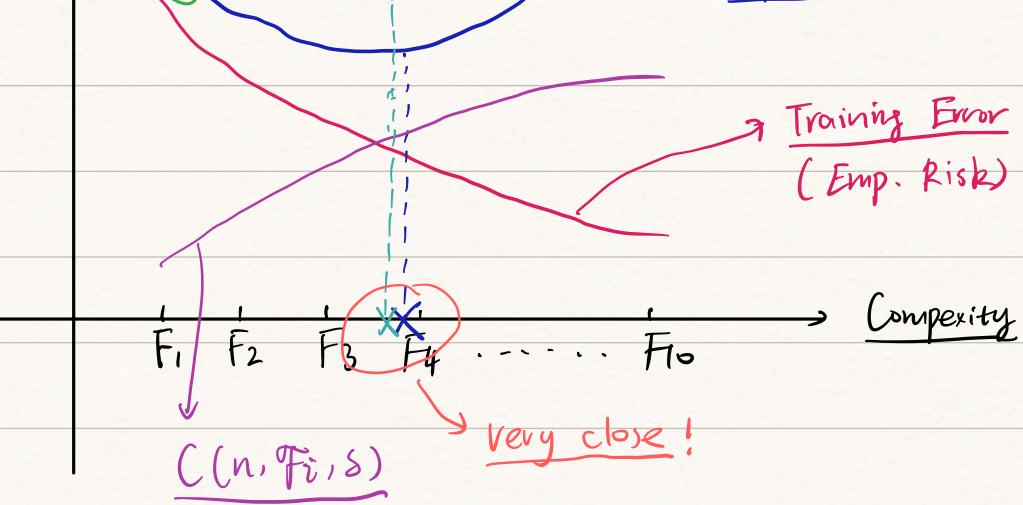
$$\forall f \in \mathcal{F}_i \quad R(f) \leq R_n(f) + C(n, \mathcal{F}_i, \delta)$$

$$C(n, \mathcal{F}_i, \delta) = \sqrt{\frac{\log |\mathcal{F}_i| + \log(\frac{2}{\delta})}{2n}}$$

$$\Rightarrow R(\hat{f}_i) \leq R_n(\hat{f}_i) + C(n, \mathcal{F}_i, \delta)$$

overfit





$$\rightarrow \text{ERM} \quad \hat{\theta}_i = \operatorname{argmin}_{\theta} \frac{1}{n} \sum \text{Loss}_0(y_t, f_i(x_t))$$

select the most complex model

$$\hat{i}^*_{\text{ERM}} = \operatorname{argmin}_{i=1, \dots, K} R_n(f_i) = \boxed{K}$$

$$\rightarrow \text{SRM} \quad i^* = \operatorname{argmin}_{i=1, \dots, K} R_n(f_i) + \underline{C(n, F_i, s)}$$

fit to model

penalize more complex models

Bayesian Score

Suppose that model F takes $x \in \mathbb{R}^d$ and maps to $y \in \mathbb{R}$

$$P(y | x, \underline{\theta}, \sigma^2) = N(y; \underline{\theta}^T x, \sigma^2)$$

Given a dataset $\{(x_t, y_t)\}_{t=1}^n \subset \mathbb{R}^d \times \mathbb{R}$. we want to estimate $\underline{\theta}$.

\downarrow Method

Maximize Likelihood!

$$\text{Likelihood} : \mathcal{L}(\underline{\theta}; \underline{x}) = \prod_{t=1}^n P(y_t | x_t, \underline{\theta}, \sigma^2)$$

$$= \prod_{t=1}^n N(y_t; \underline{\theta}^T x_t, \sigma^2)$$

$$= \prod_{t=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(y_t - \underline{\theta}^T x_t)^2}{2\sigma^2} \right\}$$

$$X = \begin{bmatrix} x_1^T & 1 \\ \vdots & \vdots \\ x_n^T & 1 \end{bmatrix}_{n \times (d+1)} \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

$$\Rightarrow \text{MLE } \hat{\underline{\theta}} = \underset{\underline{\theta}}{\operatorname{argmin}} \mathcal{L}(\underline{x}_n; \underline{\theta})$$

$$= (X^T X)^{-1} X^T \underline{y}. \quad (\text{Assume } X \text{ - full column rank})$$

Here we have point estimate of $\hat{\underline{\theta}}$, namely $(X^T X)^{-1} X^T \underline{y}$



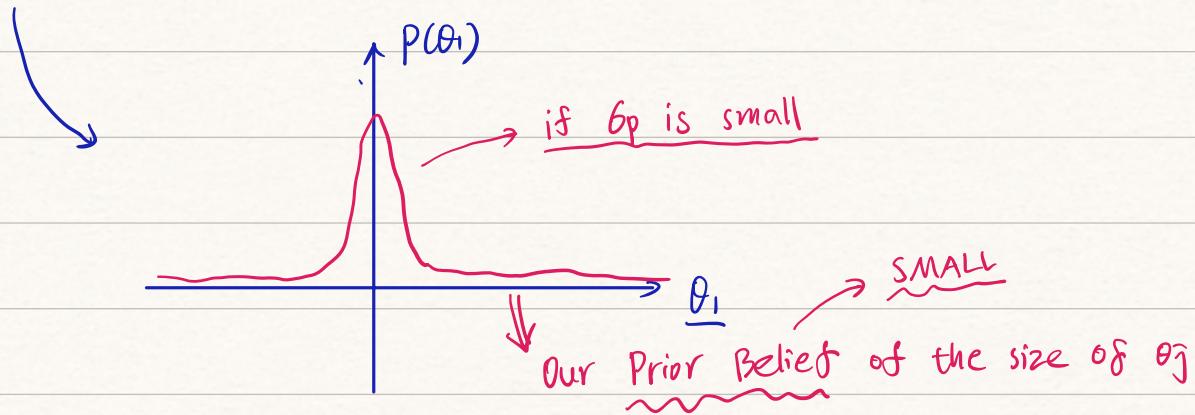
In Bayesian analysis, we want a dist. over $\underline{\theta}$ given \underline{x}_n

$$\underline{\mathcal{D}} = \{(x_t, y_t)\}_{t=1}^n \quad (x_t, y_t) \sim P$$

Method: impose a prior distribution over $\underline{\theta} \sim P(\underline{\theta}) = N(\underline{\theta}; 0, \sigma_p^2 I)$

$$P(\underline{\theta}) = \prod_{j=1}^d N(\theta_j; 0, \sigma_p^2) = \prod_{j=1}^d \frac{1}{\sqrt{2\pi\sigma_p^2}} \exp \left\{ -\frac{1}{2\sigma_p^2} \theta_j^2 \right\}.$$

Rmk: σ_p is small, we expect that each component of $\underline{\theta}$ is small



Bayesian Rule (Fundamental relation between Bayesian statistics)

$$\text{Posterior} := \frac{\text{Likelihood} \times \text{Prior}}{\text{Model Evidence}}$$

$$P(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta) P(\theta)}{P(\mathcal{D})}$$

~~Approximately~~ → Equivalent

Rmk: Model evidence \Leftrightarrow Bayesian Score (BIC is approximat-r)
 equivalent \Leftrightarrow Marginal likelihood
 (two names)

$$P(\mathcal{D}) = P(\mathcal{D} | \mathcal{F}) = \int L(\mathcal{D}; \underline{\theta}') P(\underline{\theta}') d\underline{\theta}'$$

Model class currently assuming

$$\rightarrow \hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta | \mathcal{D})$$

depend on the f class \mathcal{F}

$$= \arg \max_{\theta} \underbrace{L(\mathcal{D}; \theta)}_{\text{prior } (\text{before } \mathcal{D})} \underbrace{P(\theta)}_{\text{prior } (\text{before } \mathcal{D})}$$

after Observation

In the linear Gaussian Model

$$\begin{cases} P(\theta) = N(\theta; 0, \sigma^2_p I) \\ L(\theta; D) = \prod_{t=1}^n N(y_t; x_t^\top \theta, \sigma^2) \end{cases}$$



$$\log P(D|F) = -\frac{n}{2} \log(2\pi\sigma^2) + \frac{d}{2} \log \lambda - \frac{1}{2} \log |X^\top X + \lambda I|$$

$$\boxed{\lambda = \frac{\sigma^2}{\sigma_p^2}}$$



$$- \frac{1}{2\sigma^2} (\|y\|^2 - y^\top X(X^\top X + \lambda I)^{-1} X^\top y)$$

Actually this is very hard to get

↗ Laplace Approximation

For the linear Gaussian Model prior on θ , we can find the Model Evidence in CLOSED FORM.



Generally, it is impossible for us to calculate the Model Evidence directly!



Laplace Approximation

Linear Gaussian Model

$$P(\theta|D) = \frac{P(\theta) L(D|\theta)}{P(D)} \quad \rightarrow \text{posterior distribution}$$

$$P(\theta|D) \sim N(\theta; \mu, \Sigma)$$

$$\begin{cases} \mu = (X^T X + \lambda I)^{-1} X^T y \\ \Sigma = \sigma^2 (X^T X + \lambda I)^{-1} \end{cases}$$

Rmk: $\Sigma = \sigma^2 (X^T X + \lambda I)^{-1}$

i) depends on the original noise variance σ^2

ii) $X^T X = \sum_{t=1}^n x_t x_t^T \Rightarrow \frac{1}{n} X^T X = \frac{1}{n} \sum x_t x_t^T$

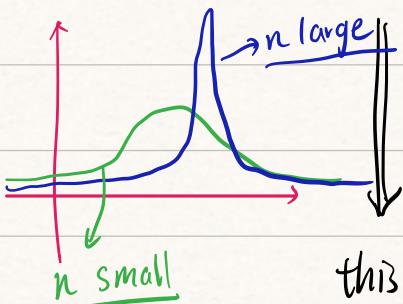
\rightarrow If $x_t \sim P(x)$ i.i.d. manner

then $\frac{1}{n} \sum x_t x_t^T \xrightarrow{P} C =$

$$\Sigma = \sigma^2 (X^T X + \lambda I)^{-1} = \sigma^2 \left(n \left(\frac{1}{n} X^T X + \frac{1}{n} \lambda I \right) \right)^{-1}$$

$$\approx \frac{\sigma^2}{n} (C + D)^{-1}$$

$$= \frac{\sigma^2}{n} C^{-1}$$

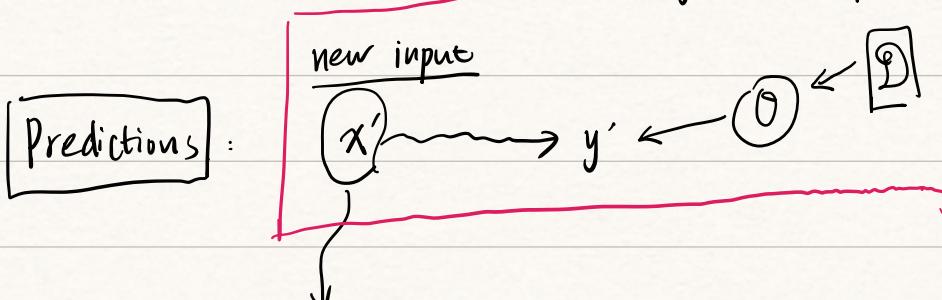


this means when we have MORE samples,

the uncertainty about our θ decays on $O(\frac{1}{n})$



We have a distribution (not just a point estimate) for $\hat{\theta}_{MAP}$



$$P(y'|x', D) = \int P(y'|\theta | x', D) d\theta$$

distribution over target

$$= \int p(y' | x', \theta) P(\theta) d\theta$$

posterior

$P(\theta | D)$

When x' is the test sample

likelihood ~~prior~~

How to select a model?



Suppose that there are 2 Competing Models!

$$f_1: P(y|x, \underline{\theta}^1, \sigma^2) = N(y; \theta^1 \phi^{(1)}(x), \sigma^2) \quad \theta \in \mathbb{R}^{d_1}$$

$$f_2: P(y|x, \underline{\theta}^2, \sigma^2) = N(y; \theta^2 \phi^{(2)}(x), \sigma^2) \quad \theta \in \mathbb{R}^{d_2}$$

$\phi^{(1)}$ ~ a linear model

$$\phi^{(1)}(x) = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$\phi^{(2)}$ ~ a quadratic model

$$\phi^{(2)}(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$f_1 \subseteq f_2$$

If can compute model evidence / Bayesian score $P(D|f_i)$

we should select the larger Model Evidence!



$$i^* = \underset{i=1, \dots, K}{\operatorname{argmax}} P(D|f_i)$$

看到 D 的可能性

Difficult to Compute

BIC

Method: Bayesian Information Criterion $\approx P(\mathcal{D} | F_i)$

Conclusion: Works well when d_i is small and $n \rightarrow \infty$

When $n \rightarrow \infty$, under some regularity conditions on the likelihood and prior terms

$$\downarrow \quad \begin{array}{c} \text{model evidence} \\ \text{log-likelihood + penalty} \end{array}$$

$$i^* = \underset{i=1,2,\dots,K}{\operatorname{arg\,max}} P(\mathcal{D} | F_i) \approx BIC(i)$$

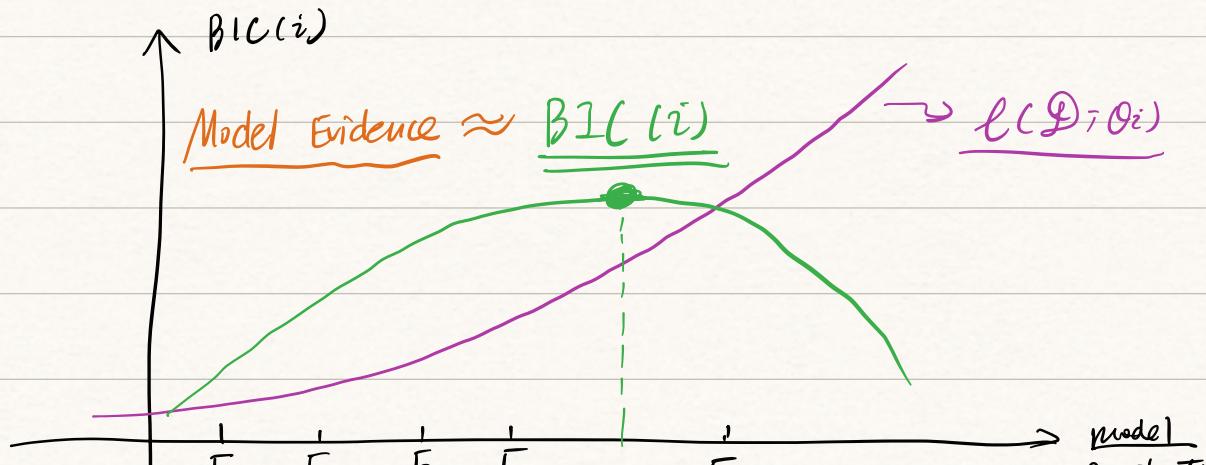
$$BIC(i) = \ell(\mathcal{D}; \hat{\theta}_i) - \frac{d_i}{2} \log n$$

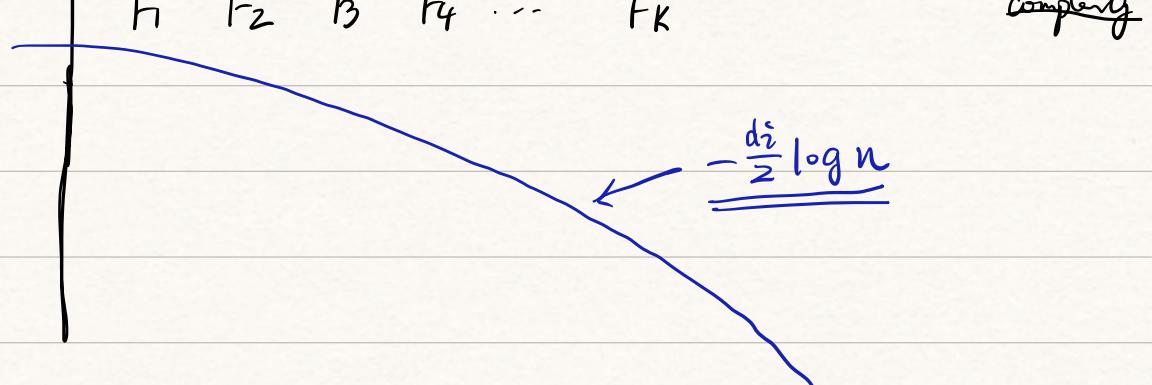
maximize log-likelihood under model F_i

$$\hat{\theta}_i = \underset{\text{log-likelihood}}{\operatorname{arg\,max}} \ell(\mathcal{D}; \theta)$$

penalization terms

Intuition: i) BIC penalizes models with a Large # of parameters ; $d_i \uparrow$, $BIC(i) \downarrow$
 ii) without the complexity penalty terms ($\frac{d_i}{2} \log n$), we will always choose the most complex model





Example : $\ell_1 = -100$ $\ell_2 = -50$

$$\ell_i = \ell(\theta; \hat{\theta}_i)$$

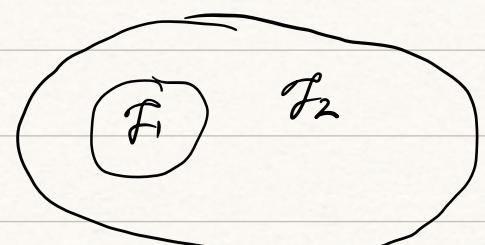
↓
 ML
 Negative $\rightarrow \log P < 0$
 ↓
 $\ell < 0$

② Linear Model) $d_1 = d+1$

Quad. Model $d_2 = \binom{d+2}{2} = \frac{1}{2}(d+1)(d+2)$

1. Why is $\tilde{\ell}_1 \leq \tilde{\ell}_2$ for $\mathcal{F}_1 \subseteq \mathcal{F}_2$?

Pf: $\tilde{\ell}_1 = \arg \max_{\theta, \mathcal{F}_1} \log L(\theta; \hat{\theta})$



$\tilde{\ell}_2 = \arg \max_{\theta, \mathcal{F}_2} \log L(\theta; \hat{\theta})$

$$\Rightarrow \tilde{\ell}_1 \leq \tilde{\ell}_2$$

2. Prefer Model 2 over Model 1 (when)

i.e. $BIC(2) > BIC(1)$

$$BIC(2) = \hat{\ell}_2 - \frac{d_2}{2} \log n = -50 - \frac{(d+1)(d+2)}{4} \log n$$

$$BIC(1) = \hat{\ell}_1 - \frac{d_1}{2} \log n = -100 - \frac{d+1}{2} \log n$$

★ { if n small \rightsquigarrow choose Model 2.
if n big \rightsquigarrow choose Model 1

#