

"Plug-in" Lecture 7
approximation

Consider predictive prob. distribution:

$$p(\bar{y}|\bar{x}, \mathcal{D}) = \int d\bar{\theta} p(\bar{y}|\bar{x}, \bar{\theta}) p(\bar{\theta}|\mathcal{D}).$$

$\int d\bar{\theta}$ hard to compute in general

If we replace $p(\bar{\theta}|\mathcal{D})$ by $\delta(\bar{\theta} - \bar{\theta}^{\text{MAP}})$,
 [or even $\delta(\bar{\theta} - \bar{\theta}^*)$]
 we obtain:

$$p(\bar{y}|\bar{x}, \mathcal{D}) \underset{==}{\approx} p(\bar{y}|\bar{x}, \bar{\theta}^{\text{MAP}}).$$

Ex. binary classification, 1D inputs
 $y \in \{0, 1\}$

$$p(y|x; \bar{\theta}) = \text{Ber}(y|G(\omega x + b)) .$$

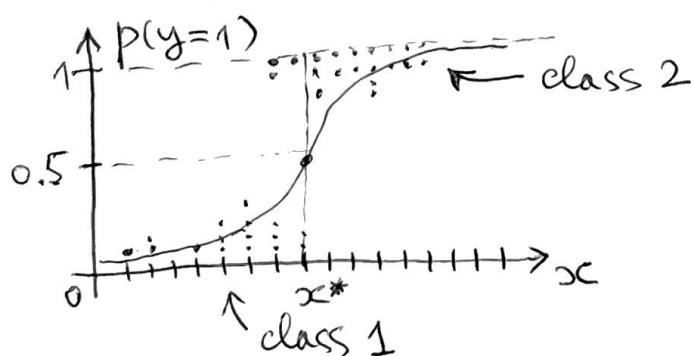
$$\bar{\theta} = \{\omega, b\}$$

In other words,

$$p(y=1|x; \bar{\theta}) = G(b + \omega x) = \frac{1}{1 + e^{-(b + \omega x)}}.$$

If we fit $\bar{\theta}$ by ML on some data

$$\mathcal{D} = \{x_n, y_n\}_{n=1}^N, \text{ we obtain:}$$



Decision boundary (DB): value of x^* s.t.

$$p(y=1|x^*; \vec{\theta}^*) = p(y=0|x^*; \vec{\theta}^*) = 0.5 .$$

Here, $G(b^* + w^*x^*) = 0.5 \Rightarrow b^* + w^*x^* = 0$, or

$$x^* = -\frac{b^*}{w^*}$$

A better approach would be to integrate over $p(\vec{\theta}|\mathcal{D})$. In practice, this can be accomplished using Monte-Carlo (MC) approximation:

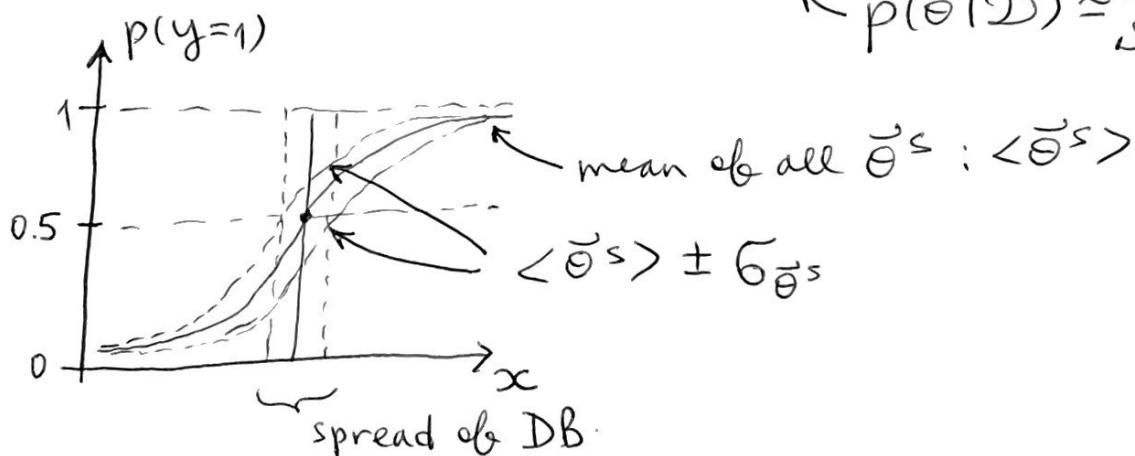
1. Draw S samples from the posterior:

$$\vec{\theta}^s \sim p(\vec{\theta}|\mathcal{D}) \quad s=1, \dots, S$$

2. Approx. pred. prob. as

$$p(y=1|x, \mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S p(y=1|x, \vec{\theta}^s)$$

$$\nwarrow p(\vec{\theta}|\mathcal{D}) \approx \frac{1}{S} \sum_{s=1}^S \delta(\vec{\theta} - \vec{\theta}^s)$$



Laplace approximation

Consider $p(\vec{\theta} | \mathcal{D}) = \frac{p(\mathcal{D} | \vec{\theta}) p(\vec{\theta})}{\underbrace{p(\mathcal{D})}_{=Z}} = \frac{p(\vec{\theta}, \mathcal{D})}{Z} =$

$$= \frac{1}{Z} e^{-\ell(\vec{\theta})}, \text{ where } \ell = -\log p(\vec{\theta}, \mathcal{D}).$$

Now, expand

$$\ell(\vec{\theta}) \approx \ell(\hat{\vec{\theta}}) + (\vec{\theta} - \hat{\vec{\theta}})^T \cdot \vec{g} + \frac{1}{2} (\vec{\theta} - \hat{\vec{\theta}})^T \cdot H \cdot (\vec{\theta} - \hat{\vec{\theta}}),$$

where

$$\left\{ \begin{array}{l} g_i = \frac{\partial \ell(\vec{\theta})}{\partial \theta_i}, \quad \text{gradient} \\ H_{ij} = \frac{\partial^2 \ell(\vec{\theta})}{\partial \theta_i \partial \theta_j} \quad \text{Hessian} \end{array} \right.$$

If we expand around the minimum of $\ell(\vec{\theta})$, we have $g_i = 0$, $\forall i$ and

$$p(\vec{\theta} | \mathcal{D}) \approx \frac{1}{Z} e^{-\ell(\hat{\vec{\theta}})} e^{-\frac{1}{2} (\vec{\theta} - \hat{\vec{\theta}})^T \cdot H \cdot (\vec{\theta} - \hat{\vec{\theta}})}$$

$$\int d\vec{\theta} p(\vec{\theta} | \mathcal{D}) = 1 \quad \text{yields}$$

$$\frac{e^{-\ell(\hat{\vec{\theta}})}}{Z} = \left[(2\pi)^{D/2} |H|^{-1/2} \right]^{-1}, \text{ or}$$

$$Z = e^{-\ell(\hat{\vec{\theta}})} (2\pi)^{D/2} |H|^{-1/2}.$$

$$\text{Finally, } p(\vec{\theta} | \mathcal{D}) = N(\vec{\theta} | \hat{\vec{\theta}}, \underline{H^{-1}}).$$

Bayesian decision theory

Classification: decide the optimal class label to predict, given an observed input \vec{x} .

Class labels: $\mathcal{Y} = \{1, \dots, C\}$

If $C=2$, we can define a zero-one loss

$$l_{01}(y, \hat{y}) = \mathbb{I}(y \neq \hat{y})$$

↑ ↑
observation model

		$\hat{y}=0$	$\hat{y}=1$
$y=0$	$\hat{y}=0$	0	1
	$\hat{y}=1$	1	0

Then the posterior expected loss

$$\begin{aligned} R(\hat{y}|\vec{x}) &= 1 \times p(\hat{y} \neq y|\vec{x}) + 0 \times p(\hat{y} = y|\vec{x}) \\ &= 1 - p(\hat{y} = y|\vec{x}). \end{aligned}$$

Thus, to minimize $R(\hat{y}|\vec{x})$ we need to maximize $p(\hat{y} = y|\vec{x})$:

$$\underset{y \in \mathcal{Y}}{\operatorname{argmax}} p(y|\vec{x}) \quad \leftarrow \begin{array}{l} \text{choose the} \\ \text{most probable} \\ \text{value of } y \\ (\text{the MAP estimate}) \end{array}$$

Now, consider the following loss matrix:

	$\hat{y} = 0$	$\hat{y} = 1$
$y = 0$	l_{00}	l_{01}
$y = 1$	l_{10}	l_{11}

If $p_0 = p(\hat{y} = 0 | \bar{x})$ and $p_1 = 1 - p_0$,

$$R(\hat{y} = 0 | \bar{x}) = l_{00} p_0 + l_{10} p_1,$$

$$R(\hat{y} = 1 | \bar{x}) = l_{01} p_0 + l_{11} p_1$$

Then we choose $\hat{y} = 0$ if

$$l_{00} p_0 + l_{10} p_1 < l_{01} p_0 + l_{11} p_1.$$

If $l_{00} = 0$, $l_{11} = 0$ (no penalty for correct decisions),

$$l_{10} p_1 < l_{01} (1 - p_1), \text{ or}$$

$$p_1 < \frac{l_{01}}{l_{10} + l_{01}}$$

l_{01} = false positive
($\hat{y} = 1$ but $y = 0$)

l_{10} = false negative
($\hat{y} = 0$ but $y = 1$)

If, in addition, $l_{10} = c l_{01}$, we pick

$$\hat{y} = 0 \text{ iff } p_1 < \frac{1}{1+c}.$$

For ex., if $c = 2 \Rightarrow p_1 < \frac{1}{3}$, or $p_0 > \frac{2}{3}$.

So, we declare $\hat{y} = 0$ if $p_0 > \frac{2}{3}$, otherwise we declare $\hat{y} = 1$.

Note that the 'old' rule (take the \hat{y} with max prob.) corresponds to $c=1$.

Consider thresholding $p(\hat{y}=1|\vec{x})$:

$p(\hat{y}=1|\vec{x}) \geq 1-\tau$, and define

$$\hat{y}_\tau(\vec{x}) = \mathbb{I}(p(\hat{y}=1|\vec{x}) \geq 1-\tau).$$

Next, for binary data $\{y_n, \vec{x}_n\}_{n=1}^N$ (^{target})
define $FP_\tau = \sum_{n=1}^N \mathbb{I}(\hat{y}_\tau(\vec{x}_n) = 1, y_n = 0)$,

and similarly $TP_\tau, TN_\tau, FN_\tau$.

		Estimate		Row sum	$\left. \begin{array}{l} \text{normalization} \\ \text{for } p(\hat{y} y) \end{array} \right\}$
		0	1		
Data	0	TN_τ	FP_τ	N	$\left. \begin{array}{l} \text{normalization} \\ \text{for } p(y \hat{y}) \end{array} \right\}$
	1	FN_τ	TP_τ	P	
Column sum		\hat{N}	\hat{P}	normalization for $p(y \hat{y})$	

		Estimate \hat{y}		' $p(\hat{y} y)$ estimates'	
		0	1		
Data	0	$\frac{TN}{N} = TNR$	$\frac{FP}{N} = FPR$	R stands for 'rate'	
	1	$\frac{FN}{P} = FNR$	$\frac{TP}{P} = TPR$		

$\left\{ \begin{array}{l} TNR + FPR = 1, \\ FNR + TPR = 1 \end{array} \right.$

-6-

Alternatively, we can focus on
' $p(y|\hat{y})$ estimates':

		Estimate \hat{y}	
		0	1
Data	0	$\frac{TN}{N} = NPV$	$\frac{FP}{\hat{P}} = FDR$
	1	$\frac{FN}{N} = FOR$	$\frac{TP}{\hat{P}} = PPV$

NPV = negative predictive value

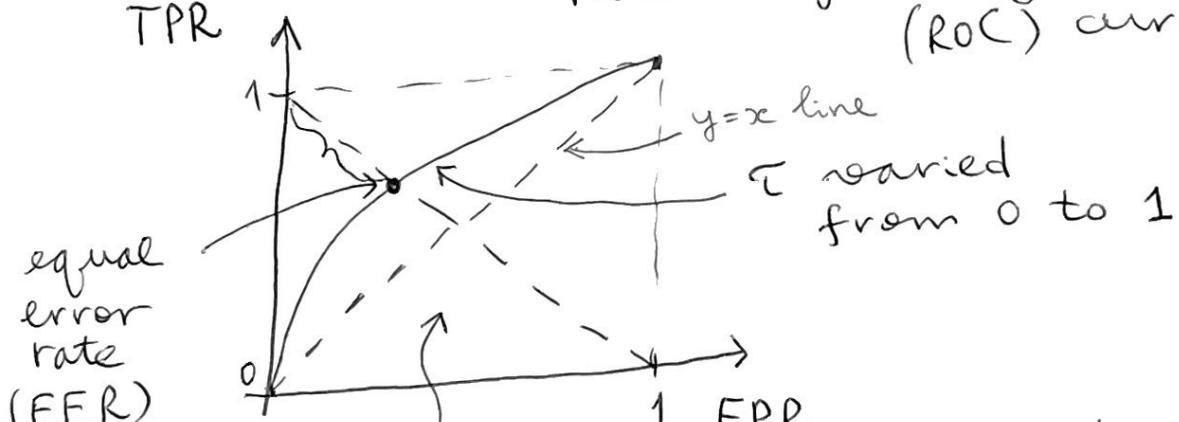
PPV = pos. pred. value, or precision

FDR = false discovery rate

FOR = false omission rate

Consider $\begin{cases} TPR_{\tau} = p(\hat{y}=1|y=1, \tau) = \frac{TP_{\tau}}{TP_{\tau} + FN_{\tau}}, \\ FPR_{\tau} = p(\hat{y}=1|y=0, \tau) = \frac{FP_{\tau}}{FP_{\tau} + TN_{\tau}}. \end{cases}$

Receiver operating characteristic (ROC) curve



area = AUC or 'ROC AUC'

-7- ROC AUC = 1.0 is the best value

EER is defined as $FPR = FNR$:

$$FPR = 1 - TPR \Rightarrow TPR = 1 - FPR.$$

$\underline{\text{EER}}$ = distance between top left corner and the EER point; 0 is the best value.

$\underline{0}$

Precision-recall curves

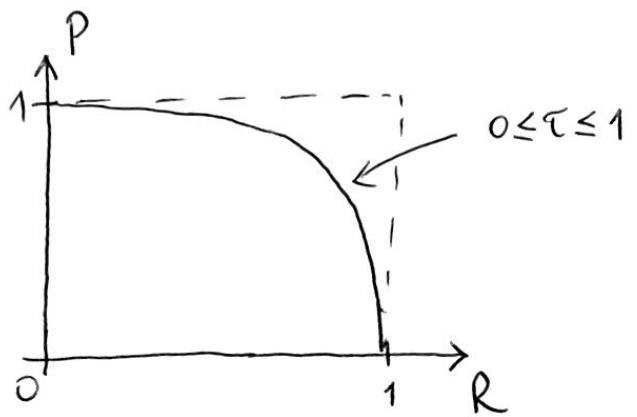
Sometimes, the notion of "negative observation" is not well-defined: e.g. in detecting objects in images, the ^{total} number of ~~object~~ patches is a patch by patch, hyperparameter, so the number of negative observations is user-defined.

In such situations, it may be useful to replace FPR with

$$\begin{aligned} \text{precision} &= PPV \equiv P(\tau) = \\ &= p(y=1 | \hat{y}=1, \tau) = \frac{TP_{\tau}}{TP_{\tau} + FP_{\tau}}. \end{aligned}$$

Likewise,

$$\begin{aligned} \text{recall} &= TPR \equiv R(\tau) = p(\hat{y}=1 | y=1, \tau) = \\ &= \frac{TP_{\tau}}{TP_{\tau} + FN_{\tau}}. \end{aligned}$$



The precision-recall curve can be summarized as:

- (1) area under the curve
(1.0 is the max)
- (2) precision of the first K entities recalled
"precision at K " score

[Precision - recall curve]

revisited

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

FP = # negative class instances
mistakenly identified as positive

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

FN = # positive class instances
mistakenly identified as negative

Thus , Precision = $\frac{\text{Correctly predicted pos.}}{\text{Total predicted pos.}}$

Recall = $\frac{\text{Correctly pred. pos.}}{\text{Truly pos. instances}} \quad (\tau = \frac{1}{2}, \text{say})$

Ex.	Document	Ground truth		Prediction	
		S	N	S	TP
Sports/Not sports	1	S		S	TP
	2	S		S	TP
	3	N		S	FP
	4	S		N	FN
	5	N		N	FN
	6	S		S	FP
	7	N		S	FP
	8	N		N	

$$\text{TP} = 2, \text{ FP} = 2, \text{ FN} = 2 \Rightarrow \left\{ \begin{array}{l} \text{Pr} = \frac{1}{2} \\ \text{Rc} = \frac{1}{2} \end{array} \right.$$

To produce the PR curve, just change the value of the threshold τ . -10-