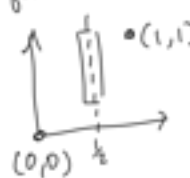First, quick review of last class's concept check



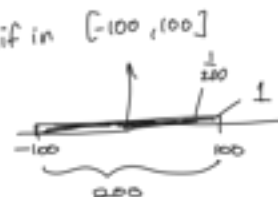$\Rightarrow$ noted small noise $\rightarrow$ parabola will go through the pts.

let's suppose: $\hat{y} = a_0 + a_1 x + a_2 x^2$

plug in $(0,0)$: $a_0 = 0$
plug in $(1,1)$: $1 + a_1 + a_2$ ; write : $a_2 = 1 - a_1$

put an improper prior on $a_2 \rightarrow$ easier $a_2$: unif in $[-100, 100]$
$p(a_2) = \frac{1}{200}$



$a_1$: unif in $[-1, 1]$   $p(a_1) = \frac{1}{2}$
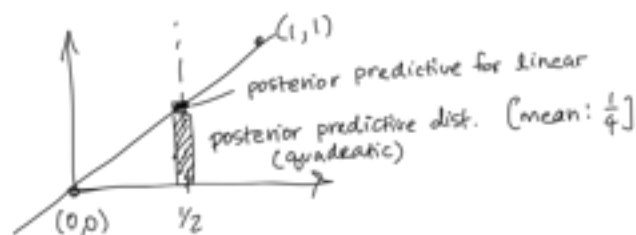
what is the posterior predictive for $x = \frac{1}{2}$

$$y(x = \frac{1}{2}) = a_1(\frac{1}{2}) + a_2(\frac{1}{4})$$
$$= a_1(\frac{1}{2}) + (1-a_1)(\frac{1}{4})$$
$$= \frac{1}{4} + \frac{1}{4}a_1 \quad \} \text{ linear ; } \underline{\text{uniform}} \text{ prior over } a_1$$

$a_1 = -1 : \hat{y} = 0$
$a_1 = 1 : \hat{y} = \frac{1}{2}$



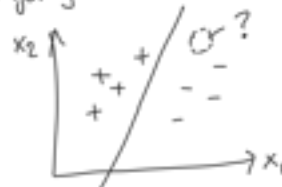posterior predictive for linear
posterior predictive dist. [mean: $\frac{1}{4}$]
(quadratic)

if we had a linear model?

## Onto classification

$\hat{y}$ is going to be a categorical variable.



Conceptually same framework as regression...
$\rightarrow$ we need decide on a model class? (what kind of separators?)
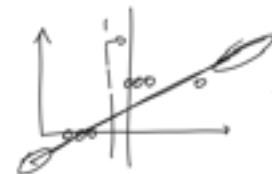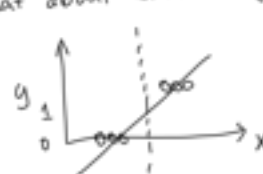$\rightarrow$ we need decide on a loss function?

Notation: $\boxed{\hat{y} \in \underbrace{C_1 \ldots C_k}_{\substack{\text{discrete/}\\\text{countable.}}}}$   if $y$ is binary, $y = -1/1$
or
$y = 0/1$

if $y$ is k-any ; $[0 \quad 0 \quad 1 \quad 0 \quad 0]$
$\underbrace{\quad\quad\quad\quad}_{\text{"one hot" encoding}}$

Can we just use our regression tools?
$\rightarrow$ KNN still works (nonparametric regression $\rightarrow$ classification)
(avg)                              (majority vote)

$\rightarrow$ what about linear regression?



$\Rightarrow$ can't just apply linear reg to the 0|1 encoding of classes.

## Linear Classification

$$\hat{y} = \text{sign}(w^T x + w_0) = \text{sign}(w^T x) \qquad \text{here } \hat{y} \in \{1,1\}$$

$$\text{if } w^T x < 0 \ , \quad -1 = \hat{y}$$
$$w^T x > 0 \ , \quad 1 = \hat{y}$$

$w^T x + w_0 = 0$ is the $\underline{\text{decision boundary}}$

in 2D case: $0 = w_0 + w_1 x_1 + w_2 x_2$

$$x_2 = \frac{-w_1 x_1}{w_2} + \frac{-w_0}{w_2} \quad \Big] \text{ clearly a line}$$

$\underbrace{\qquad}_{\text{slope}}$ $\underbrace{\qquad}_{\text{offset}}$


offset $\frac{-w_0}{w_2}$   slope $\frac{-w_1}{w_2}$

generalize to multiple dims (D>2): consider a vector connecting 2 points on the boundary:



we know: $\begin{array}{l} w^T x_1 + w_0 = 0 \\ w^T x_2 + w_0 = 0 \end{array} \Big] \begin{array}{l}\text{by}\\\text{definition}\end{array}$

consider $s = x_2 - x_1$ lies on the boundary

$$w^T s = w^T x_2 - w^T x_1$$
$$= \underbrace{w^T x_2 + w_0}_{0} \underbrace{- w_0 - w^T x_1}_{0}$$
$$= 0$$

implication is that $w$ is orthogonal to the boundary

$\underline{\text{Now}}$: let's introduce a loss function. [esp. important if we cannot perfectly separate the data]:



- $l_{0/1}(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{else} \end{cases}$

$$\underline{L(w) = \sum_n l_{0/1}\left(-y_n\left(w^T x_n + w_0\right)\right)}$$

$\underbrace{\phantom{}}_{\substack{\text{true}\\\text{value}\\(\pm 1)}}$ $\underbrace{\phantom{}}_{\substack{\text{predicted}\\\text{value, before}\\\text{sgn function}}}$

this is positive $\underline{\underline{\text{IF}}}$ sign of
$y_n \neq$ sign of $w^T x_n + w_0$

$\Rightarrow$ uninformative gradient!

$= \#$ of mis-classified examples.

- alternative : hinge loss / rectified linear function

$$l_{hinge} = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{else} \end{cases} \qquad \max(0, z)$$

$$L(w) = \sum_n l_{hinge}\left(-y_n\left(w^T x + w_0\right)\right) \quad \xrightarrow{\quad} \text{ASIDE: use } \underline{\text{any}} \ f(x, w) \text{ here}$$

$$= \sum_{\substack{n \text{ that} \\ \text{are incorrect}}} -\left(w^T x_n + w_0\right) y_n$$

We can take gradients! $\frac{\partial}{\partial w} L(w) = -\underbrace{\sum y_n x_n}_{\substack{\text{bad} \\ y\text{'s}}}$

Interesting connection: old algorithm, 1958, Perceptron Algorithm $\left(\begin{array}{l}\text{for}\\\text{linear}\\\text{classification}\end{array}\right)$

- loop through the data
  - for each mis-classified pt. if $\hat{y}_n \neq y_n$,

$$w \leftarrow w + \underline{y_n x_n}$$

until no errors; will converge if a perfect separator exists.

$\underline{\text{New idea}}$: stochastic gradient descent, use a $\underline{\text{minibatch}}$ of the data when computing the derivative.

relevant: big data.

properties re convergence
⮡ look at perceptron: that's just SGD w/ minibatch size=1

Just briefly: one more loss function: if $w$ projects $x$ onto 1D, $(w^T x)$
suppose we want some kind of "separation" or "clustering" of the
different classes...



(Fisher's Discriminant)

define empirical means & variances:
$m_1 = \frac{1}{N_1} \sum x_n$ s.t. $y_n = C_1$
$m_2 = \frac{1}{N_2} \sum x_n$ s.t. $y_n = C_2$
$S_1 = \frac{1}{N_1} \sum (x_n - m_1)(x_n - m_1)^T$
$S_2 = \frac{1}{N_2} \sum (x_n - m_2)^T (x_n - m_2)$

then: $w^T x$ : $m_1' = w^T m_1$   $V_1 = w^T S_1 w$
              $m_2' = w^T m_2$   $V_2 = w^T S_2 w$

min $\mathcal{L}(w) = \dfrac{-(m_1' - m_2')^2}{V_1 + V_2}$   } put means far apart
                                                    } make variances small

we can optimize this loss in closed form: $w \propto (S_1 + S_2)^{-1}(m_1 - m_2)$

___

Take a step back & discuss additional metrics

Errors:

| $y$ | $\hat{y}$ | |
|---|---|---|
| 1 | 1 | True positive |
| 0 | 1 | False positive |
| 0 | 0 | True Neg. |
| 1 | 0 | False neg. |

⇒ option: directly put costs on each error.

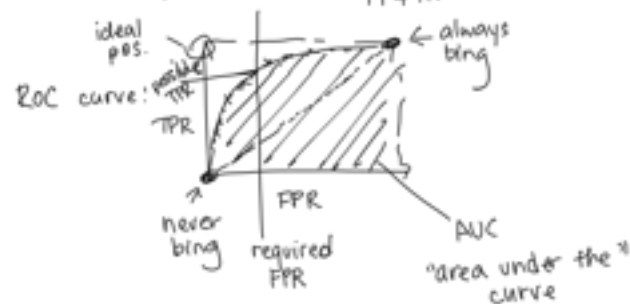now, we can combine these diff. errors to get various quantities:

accuracy: $\dfrac{TP + TN}{TP + TN + FP + FN}$

precision: $\dfrac{TP}{TP + FP}$

recall: $\dfrac{TP}{TP + FN}$  } all positives

true pos rate: $\dfrac{TP}{TP + FN}$

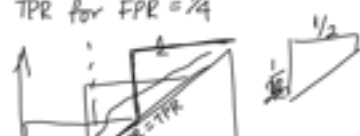false pos rate: $\dfrac{FP}{FP + TN}$



ROC curve:
ideal pos.
← always bing
TPR
never bing
FPR
required FPR
AUC
"area under the" curve

$w^T x$ is a number;
we can sweep over $w_0$
to get diff. classifiers...

may also see



prec
recall

Concept Check w/ AUCs:
suppose you have a model w/ AUC $= \frac{1}{2}\left(\frac{1}{2} + \frac{1}{8}\right) = \frac{5}{8} = .625$ (low AUC)
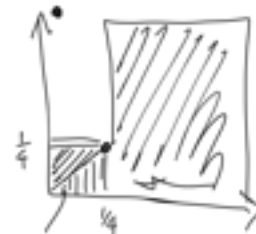
1. What are possible TPR for FPR $= \frac{1}{4}$



$TPR = \frac{1}{4}$   $FPR = \frac{1}{4}$

2. How does this change if AUC was bigger? What's the highest AUC you can have and still TPR = ¼, FPR = ¼?



$$\frac{1}{16} + \frac{3}{4} = \frac{10}{16} = .8125$$

Classification