$X_2$

$f(\ )$

Inductive Bias
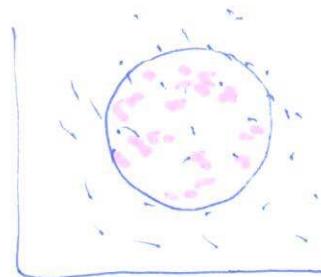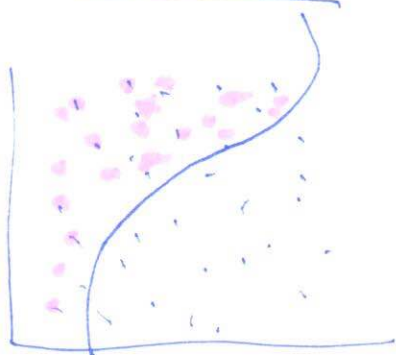
$f \rightarrow$ linear

$X_1$

Linear Boundary
Discriminating Surface
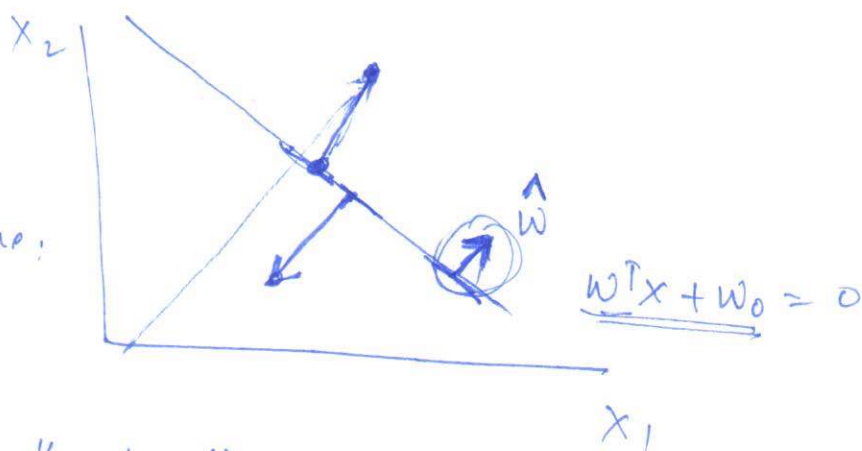
$\rightarrow$ Linearly Separable

For a point on the line:

$$w^T x + w_0 = 0$$

For a point "above" the line:

$$w^T x + w_0 > 0$$

For a point "below" the line:

$$w^T x + w_0 < 0$$

$$w^T x + w_0 = 0$$

Learning $(W, w_0)$

Generalization (True) Risk   vs.   Empirical Risk.

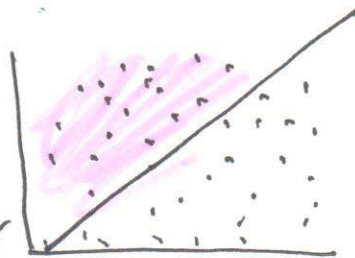loss on future (unseen) data → True Risk

We do not have data to measure true risk.

We can measure risk on training data

Empirical risk.

All data

Train data

Overfitting ←

Train data

Assume
$x_i$ do not
have bias
term

| | | $w$, $w_0$ |
|---|---|---|
| $x_1$ | $y_1$ | $\text{sign}(w^T x_1 + w_0)$ |
| $x_2$ | $y_2$ | $\text{sign}(w^T x_2 + w_0)$ |
| ⋮ | ⋮ | ⋮ |
| $x_n$ | $y_n$ | $\text{sign}(w^T x_n + w_0)$ |

$x_i$     $y_i$     If $(w^T x_i + w_0) \geq 0$    $\tilde{y}_i = +1$

If $(w^T x_i + w_0) < 0$    $\tilde{y}_i = -1$

If $y_i = +1$    and $\tilde{y}_i = +1$    ⟹ No mistake

$y_i = -1$    and $\tilde{y}_i = -1$    ⟹ no mistake

$y_i = +1$    and $\tilde{y}_i = -1$ ⟹ Mistake

$+1$    ⟹ Mistake

If      $y_i \tilde{y}_i > 0$          $\Rightarrow$ No mistake

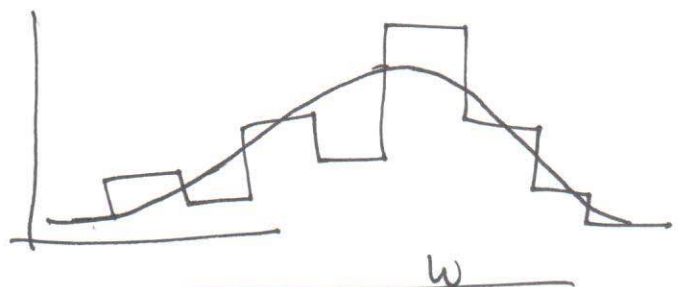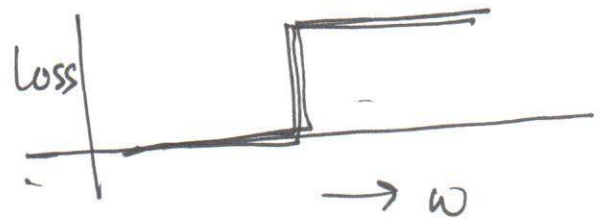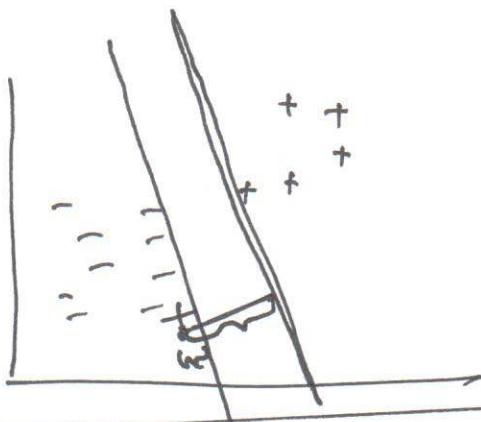       $y_i \tilde{y}_i < 0$          $\Rightarrow$ Mistake.

       $y_i (w^T x_i + w_0) > 0$        $\Rightarrow$ No mistake

       $y_i (w^T x_i + w_0) < 0$        $\Rightarrow$ Mistake.

$$\sum_{i=1}^{n} \mathbb{I}\left( y_i (w^T x_i + w_0) < 0 \right) \longrightarrow \quad 0\text{-}1 \ \ loss$$

Find $\ \underline{w, w_0}\ $ that minimizes  $0\text{-}1$ loss



let us assume that we have only one training
data point.   $x_i$    $y_i = +1$
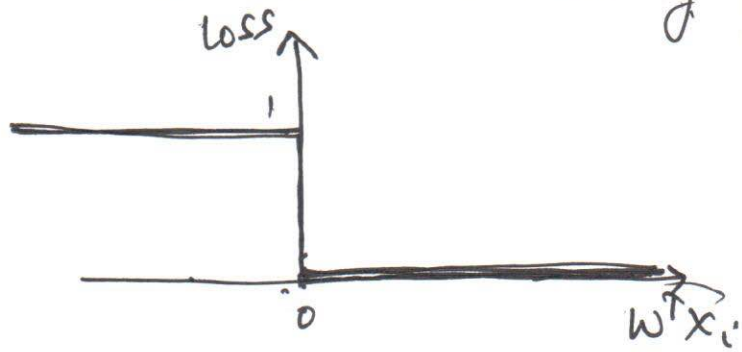
$\underbrace{(w^T x_i + w_0)}_{\theta_i}$     if $\underline{\theta_i > 0}$

$(y_i - w^T x_i)^2 = (1 - \theta_i)^2$

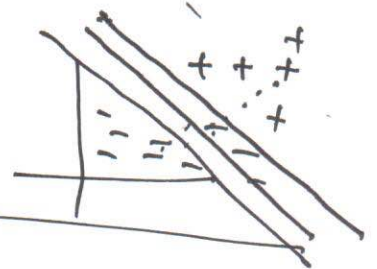$(1 + exp(-y_i w^T x_i))$

let us assume that we have 1 training point.

let $y_i = +1$

$\xrightarrow{\hspace{1cm}} \underline{w^T x_i}$

Contains
the bias!



$$(y_i - w^T x_i)^2$$

$$SL - (1 - \underline{w^T x_i})^2$$

Perception



$$\log(1 + \exp(-y_i w^T x_i))$$

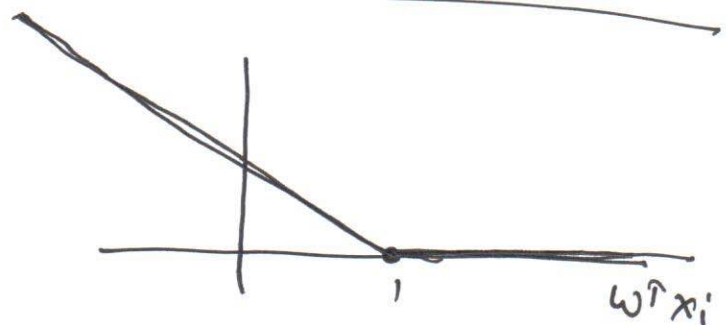Log Loss. $\log(1 + \exp(-\underline{w^T x_i}))$

logistic Regression



$$\max(0, 1 - y_i w^T x_i)$$

$$\max(0, 1 - \underline{w^T x_i})$$

Hinge loss.

Support Vector Machine

Sigmoid. $\sigma(a) = \dfrac{1}{1+\exp(-a)} = \dfrac{1}{1+e^{-a}}$

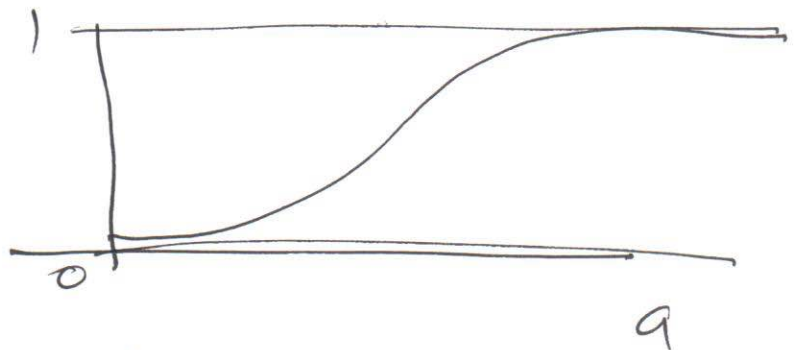$\dfrac{d}{da}\sigma(a) = \dfrac{\times 1}{(1+e^{-a})^2} \times e^{-a} = \dfrac{e^{-a}}{(1+e^{-a})^2}$

$\qquad = \dfrac{e^{-a}}{(1+e^{-a})} \ast \dfrac{1}{(1+e^{-a})}$

$\qquad = \left(1 - \dfrac{1}{1+e^{-a}}\right)\left(\dfrac{1}{1+e^{-a}}\right)$

$\qquad = \sigma(a)(1 - \sigma(a))$

$x_i$

$\sigma(w^T x_i) = \boxed{\dfrac{1}{1+\exp(-w^T x_i)}}$



$P(y_i = +1) = \sigma(w^T x_i)$

$P(y_i = -1) = 1 - \sigma(w^T x_i)$

| $x_1$ | $y_1$ | $p_1$ |
| $x_2$ | $y_2$ | $p_2$ |
| $x_3$ | $y_3$ | $p_3$ |
| $\vdots$ | $\vdots$ | $\vdots$ |
| $x_n$ | $y_n$ | $p_n$ |

$p_i = \dfrac{1}{1+\exp(-y_i w^T x_i)}$

$p_1 \times p_2 \times \cdots \cdots p_n$

$J(w) = -\sum \log p_i$

$\qquad = \sum \log\left[1 + \exp(-y_i w^T x_i)\right]$

$$J(w) = \frac{1}{n} \sum_{i=1}^{n} \log\left(1 + \exp(-w^T x_i)\right)$$

$$\nabla J(w) = \frac{1}{n} \sum_{i=1}^{n} \frac{d}{dw} \log\left(1 + \exp(-w^T x_i)\right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \frac{1}{(1 + \exp(-w^T x_i))} \exp(-w^T x_i)(-x_i)$$

$$= -\frac{1}{n} \sum_{i=1}^{n} \left[\frac{\exp(-w^T x_i)}{1 + \exp(-w^T x_i)}\right] x_i$$

Set $\nabla J(w) = 0$

and solve for $w$

─────────────────────────

Gradient descent.

$w \leftarrow w_{init}$.

until converged:

$$w \leftarrow w - \eta \nabla J(w)$$

too sensitive to $\eta$

─────────────────────────

Newton's Method.

$w \leftarrow w_{init}$

until converged

$$w \leftarrow w - \eta H^{-1} \nabla J(w)$$

Hessian Matrix