# UNIT 5 CLASSIFICATION

- predict a underline{categorical} label $y$ from $\underline{x}$ given training data $\{(\underline{x}_i, y_i)\}_{i=1}^{m}$

- binary $y$:



$y=1$ or $y \neq 1$

decision region

linear $\longrightarrow$

hyperplane decision boundary

- binary linear classification: 1) compute a "score" $z = b + \underline{w}^T \underline{x} \in \mathbb{R}$

  2) threshold: $\begin{cases} z \geq t & \Rightarrow \hat{y} = 1 \\ z < t & \Rightarrow \hat{y} \neq 1 \end{cases}$

  <span style="color:red">choose $t=0$ to maximize classification accuracy</span>

- training via logistic model

  - model: $\Pr\{y=1 \mid z\} = \dfrac{1}{1 + e^{-z}}$



"sigmoid"

  - ML estimation: $(b, \underline{w})_{ML} = \underset{(b, \underline{w})}{\arg\max} \prod_{i=1}^{n} p(y_i \mid z_i)$ for $z_i = b + \underline{w}^T \underline{x}_i$

    $= \underset{(b, \underline{w})}{\arg\min} \sum_{i=1}^{n} \left[ \ln(1 + e^{z_i}) - y_i z_i \right]$ $\longleftarrow$ <span style="color:red">for $y_i \in \{0,1\}$ binary cross entropy (BCE) loss</span>

- $K > 2$ classes

  - linear classifier: 1) $z_k = b_k + \underline{w}_k^T \underline{x} \in \mathbb{R}$ for $k = 1..K$

    2) $\hat{y} = \underset{k}{\arg\max} \ z_k$

  - training: A) one-vs-rest (OVR): for each class $k = 1..K$, train $(b_k, \underline{w}_k)$ using <span style="color:blue">binary logistic regression</span>

    B) multinomial logistic regression (MLR):

    - model: $\Pr\{y=k \mid \underline{z}\} = \dfrac{e^{z_k}}{\sum_{\ell=1}^{K} e^{z_\ell}}$ where $z_k = b_k + \underline{w}_k^T \underline{x}$

    - ML estimation: $\underset{(\underline{b}, \underline{W})}{\arg\max} \prod_{i=1}^{n} \Pr\{y_i = k \mid \underline{z}_i\}$ when $\underline{z}_i = \underline{b} + \underline{W} \underline{x}_i$
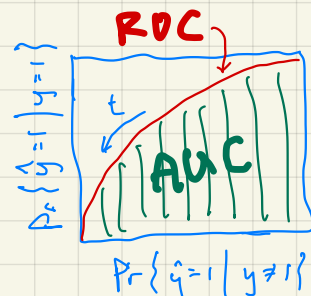
- Performance metrics (binary)

  - accuracy: $\Pr\{\hat{y} = y\}$

  - precision: $\Pr\{y = 1 \mid \hat{y} = 1\}$

  - recall: $\Pr\{\hat{y} = 1 \mid y = 1\}$

  - others...

<span style="color:red">by choosing threshold "t" we can trade off between various metrics, like precision vs recall</span>



ROC

AUC

$\Pr\{\hat{y} = 1 \mid y \neq 1\}$

<span style="color:green">AUC is a threshold-independent metric</span>

## UNIT 6 — OPTIMIZATION

• gradient: say $\underline{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$ a-d $J(\underline{w}) \in \mathbb{R}$

then $\nabla J(\underline{w}) = \begin{bmatrix} \frac{\partial J(\underline{w})}{\partial w_1} \\ \vdots \\ \frac{\partial J}{\partial w_d}(\underline{w}) \end{bmatrix}$   *same shape as $\underline{w}$*   *also can do for matrices, tensors, scalars*

   tells us the <u>slope</u> a-d <u>direction</u> of maximum increase of $J$ at $\underline{w}$

• gradient descent (GD) optimization of a smooth cost $J(\cdot)$:

   Initialize $\underline{w}$ and $\underline{w}^0$ and iterate for $t = 1, 2, 3, \dots$

   $\underline{w}^{t+1} = \underline{w}^t - \alpha_t \nabla J(\underline{w}^t)$
   *$t$ small positive stepsize*

• analysis: for sufficiently small $\alpha_t$, we proved $J(\underline{w}^{t+1}) \le J(\underline{w}^t)$ for all $t$

   via $J(\underline{w}^{t+1}) = J(\underline{w}^t) \underbrace{- \alpha_t \|\nabla J(\underline{w}^t)\|^2 + \alpha_t^2 O(\|\nabla J(\underline{w}^t)\|^2)}_{\le 0 \text{ if } \alpha_t \text{ is small enough}}$
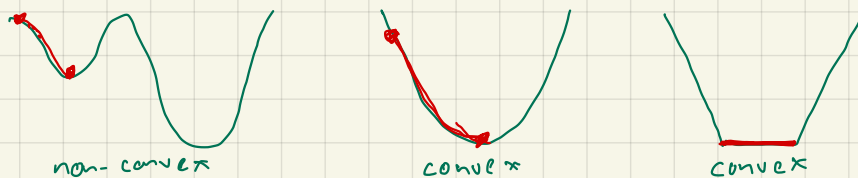
- Armijo stepsize adaptation: At each iteration $t$, perform GD update, then
   - check if $J(\underline{w}^{t+1}) \le J(\underline{w}^t) - \frac{1}{2}\alpha_t \|\nabla J(\underline{w}^t)\|^2$   *adjustable*
   - if not, set $\alpha_t \leftarrow \frac{1}{2}\alpha_t$ and try $t$'th iteration again
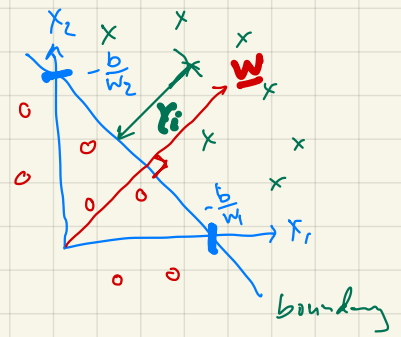   - if yes, set $\alpha_t \leftarrow 2\alpha_t$ and accept the iteration and increment $t$

- If $J(\cdot)$ is convex, then all local minimizers are global minimizers



   non-convex        convex        convex

   examples: RSS, logistic regression (and L2 & L1 regularized versions)
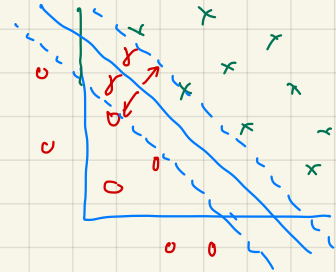   hinge loss (for SVC)

- binary linear classification

$$\gamma_i = \frac{b + \underline{w}^T x_i}{\|\underline{w}\|}$$

signed distance from $x_i$ to boundary

- "linearly separable" data means that there exists some $(b, \underline{w})$ such that

$$y_i(b + \underline{w}^T x_i) \geq 1 \quad \forall i$$

$$\iff y_i \gamma_i \geq \frac{1}{\|\underline{w}\|} \triangleq \gamma \quad \forall i \quad \longleftarrow \text{"margin"}$$

- hard margin classifier

$$\arg\min_{b, \underline{w}} \frac{1}{2}\|\underline{w}\|^2 \quad \text{s.t.} \quad y_i(b + \underline{w}^T x_i) \geq 1 \quad \forall i$$

↖ only feasible when data is linearly separable

- soft margin or support vector classifier (SVC)

$$\arg\min_{b, \underline{w}} \left\{ C \sum_{i=1}^{n} \max\{0, 1 - y_i z_i\} + \frac{1}{2}\|\underline{w}\|^2 \right\}$$

$\underbrace{\quad\quad}_{\text{hinge loss}}$  $z_i = b + \underline{w}^T x_i$

( bias-variance tradeoff )

↳ as C increases, the margin $\gamma$ decreases

the solution $(b_*, \underline{w}_*)$ depends only on the support vectors

$$\left\{ x_i : y_i(b_* + \underline{w}_*^T x_i) \leq 1 \right\}$$

- the SVM is SVC with a (usually nonlinear) transformation $\underline{x} \to \underline{\phi}(\underline{x})$
  But the SVM avoids directly computing (or even defining) $\phi(\cdot)$
  Instead it computes the kernel $K(x_i, x_j) \triangleq \underline{\phi}^T(x_i)\underline{\phi}(x_j)$
  using the "dual" form of the SVC optimization problem
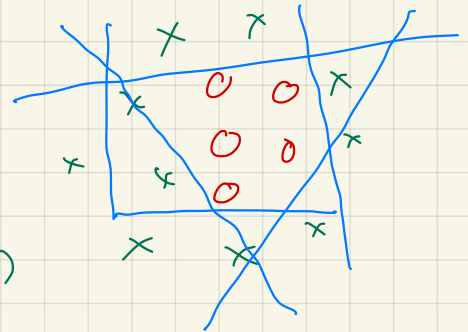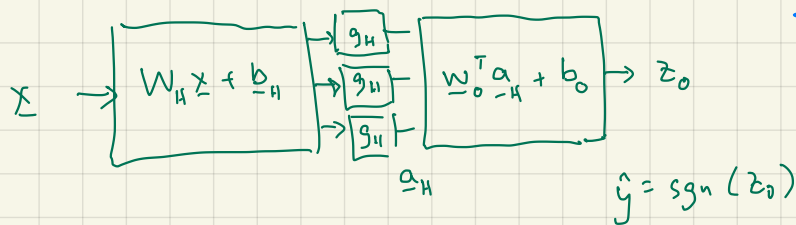
  The most popular kernel is RBF:
  $$K(x_i, x_j) = \exp\left(-\alpha\|x_i - x_j\|^2\right)$$
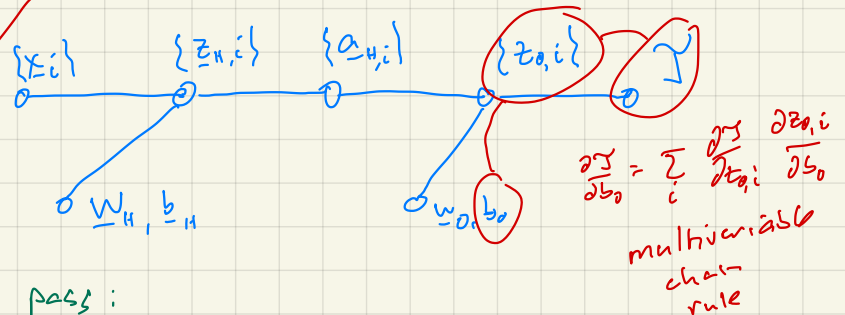  ↳ tunable $> 0$

- Can we learn a feature transformation $\phi(\cdot)$ from the training data?

- Yes, we will use the form $\left[\underline{\phi(\underline{x})}\right]_\ell = g_H\left(b_{H,\ell} + \underline{w}_{H,\ell}^T \underline{x}\right)$ for $\ell = 1 \ldots d_H$

  and then do linear classification / regression

$$\underline{x} \rightarrow \boxed{W_H \underline{x} + \underline{b}_H} \rightarrow \boxed{g_H} \rightarrow \boxed{\underline{w}_o^T \underline{a}_H + b_o} \rightarrow z_o$$

  $\underline{a}_H$

  $\hat{y} = \text{sgn}(z_o)$

- Train $\underline{\Theta} = \left\{ \underline{W}_H, \underline{b}_H, \underline{w}_o, b_o \right\}$

  to minimize $\mathcal{J}(\underline{\Theta}) \overset{\triangle}{=} \sum_{i=1}^n J(\underline{\Theta}, x_i, y_i)$

  via GD: $\underline{\Theta}^{t+1} = \underline{\Theta}^t - \alpha_t \nabla\mathcal{J}(\underline{\Theta}^t)$   or split into mini-batches

- To compute gradients

  Computation graph

  $\{x_i\}$   $\{z_{H,i}\}$   $\{a_{H,i}\}$   $\{z_{o,i}\}$   $\mathcal{J}$

  $\underline{W}_H, \underline{b}_H$   $\underline{w}_o, b_o$

  $\dfrac{\partial \mathcal{J}}{\partial b_o} = \sum_i \dfrac{\partial \mathcal{J}}{\partial z_{o,i}} \dfrac{\partial z_{o,i}}{\partial b_o}$

  multivariable chain rule

  Backpropagation:   1) forward pass:

  compute $\{z_{H,i}\}, \{a_{H,i}\}, \{z_{o,i}\}, \mathcal{J}$

  2) backward pass

  compute: $\left\{\dfrac{\partial \mathcal{J}}{\partial z_{o,i}}\right\}, \left\{\dfrac{\partial \mathcal{J}}{\partial a_{H,i}}\right\}, \left\{\dfrac{\partial \mathcal{J}}{\partial z_{H,i}}\right\}$

  $\dfrac{\partial \mathcal{J}}{\partial b_o}, \dfrac{\partial \mathcal{J}}{\partial \underline{w}_o}, \dfrac{\partial \mathcal{J}}{\partial \underline{b}_H}, \dfrac{\partial \mathcal{J}}{\partial \underline{W}_H}$