SVM objective      <span style="color:red">margin</span>      <span style="color:red">Penalty/Violations</span>

$$\min_{\substack{h \\ \vec{w}, b}} \quad \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{n} \xi_i^k$$

$$\text{s.t.} \quad y_i\left(w^T x_i + b\right) \geq 1 - \xi_i$$



<span style="color:red">ok</span>

<span style="color:red">$x_i$</span>

<span style="color:red">$\xi_i > 1$</span>

$\times$  $\times$

<span style="color:red">$0 < \xi_i \leq 1$</span>

$\xi_i = 0$

## Dual Objective

$\alpha_i$ : one for each point

$$0 \leq \alpha_i \leq C$$

$$\max_{\alpha_i} \left\{ \sum_i \alpha_i - \sum_i \sum_j \alpha_i \alpha_j \, y_i y_j \, K(x_i x_j) \right\}$$

$$\vec{w} = \sum_{\alpha_i > 0} \alpha_i y_i \, \phi(x_i)$$

linear

$$\vec{w} = \sum_{\alpha_i > 0} \alpha_i y_i x_i$$

we __Cannot__/donot have
access to $\vec{w}$
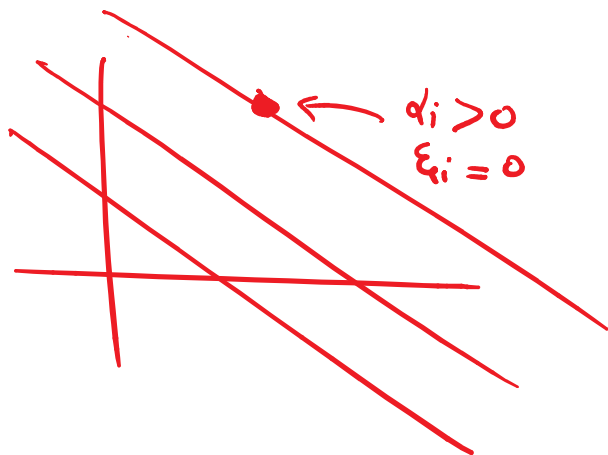
bias:
b?

$$\alpha_i \left( y_i \left( w^T x_i + b \right) - 1 + \xi_i \right) = 0$$

for support vectors
$C > \alpha_i > 0$ $\implies$ $y_i \left( w^T x_i + b \right) - 1 + \xi_i = 0$

$\hookrightarrow \xi_i = 0$



$\alpha_i > 0$
$\xi_i = 0$

$$y_i \left( w^T x_i + b \right) = 1$$

In kernel space

$$y_i \left( w^T \phi(x_i) + b \right) = 1$$

$$y_i b = 1 - y_i w^T \phi(x_i)$$

$w = \sum\limits_{\alpha_j > 0} \alpha_j y_j \phi(x_j)$

$$b = \frac{1}{y_i} - w^T \phi(x_i)$$

$$b = \frac{1}{y_i} - \left( \left( \sum \alpha_j y_j \phi(x_j) \right)^T \phi(x_i) \right)$$

$C > \alpha_i > 0$

$$b = \frac{1}{y_i} - \sum\limits_{\alpha_j > 0} \alpha_j y_j K(x_j, x_i)$$

Original
Support
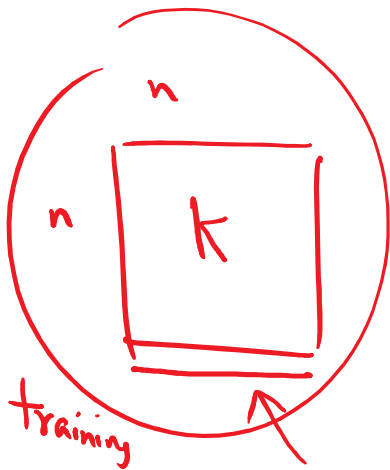
Original
Support
vectors

$$h(\vec{z}) = \underline{w^T \phi(z)} + b$$

$$= \sum_{\alpha_i > 0} \alpha_i y_i k(x_i, z) + b$$

Solved for

$z$ is a new test case



n

$k$

n

training

$k(x_i, z)$

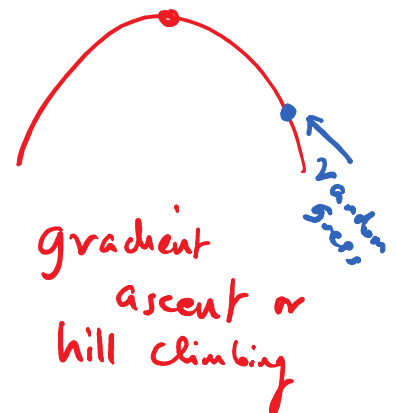$$SVM(\vec{z}) = \text{Sign}\left( h(\vec{z}) \right) - $$

$z$

---

$$\max_{\alpha_i} \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j k(x_i, x_j)$$

$$0 < \alpha_i \leq C \quad \longleftarrow \quad \text{Constraint}$$

$$\vec{\alpha} = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_n \end{pmatrix}$$
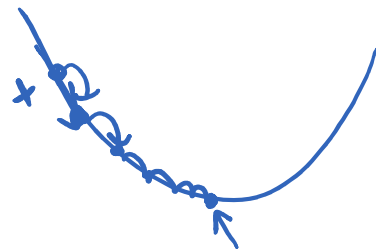
n - values to solve for

gradient ascent or hill climbing

gradient ascent

start at some solution

$\vec{\alpha}_0$, compute the gradient

$\alpha_0$, compute the gradient

derivative at $\vec{\alpha}_0$
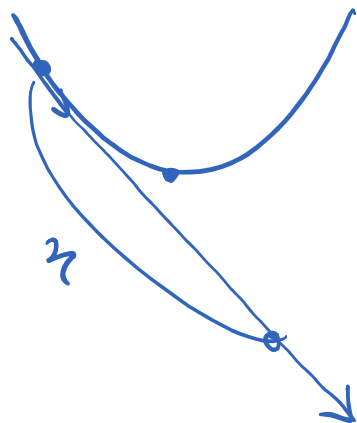
$$\nabla_{\alpha_0}$$

$$\vec{\alpha}_{i+1} = \vec{\alpha}_i + \eta \nabla_{\vec{\alpha}_i}$$

Step Size

Steepest Improvement

Sum dual objective

Convex quadratic program

$\Rightarrow$ Unique global optimal solution

$$\underline{\alpha_1 + \alpha_2 + \ldots + \boxed{\alpha_k} + \ldots}$$

$$J = \boxed{\sum \alpha_i} - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \, y_i y_j \, k(x_i x_j)$$

$$\nabla_k = \frac{\partial J}{\partial \alpha_k} = \boxed{1 - y_k \left( \sum_{i=1}^{n} \alpha_i y_i k(x_i x_k) \right)}$$

Compute gradient for k-th element of $\vec{\alpha}$

$x_k$

$$\vec{\alpha}_0 = (0, 0, \ldots, 0)$$

$$\vec{\alpha}_0 = (0, 0, \ldots, 0)$$

repeat until convergence
$\alpha_{prev} = \alpha$
for $k = 1$ to $n$ ← randomize choice of $k$

old    $\alpha_{new}$

   update the $k$-th component

$$\alpha_k = \underset{old}{\alpha_k} + \underset{\underset{step\ size}{\uparrow}}{\eta} \cdot \underset{gradient}{\nabla_k}$$

$$\alpha_k = \alpha_k + \eta \cdot \left( 1 - y_k \left( \sum_{i=1}^{n} \alpha_i y_i \ k(x_i \cdot x_k) \right) \right)$$

Clipping: $0 < \alpha_k \le C$ ← ensure

1) Approach 1: batch approach

    New values are not used until the next iteration

2) SGD: Stochastic Gradient Descent

    Using new updated values the moment they become available

    Much faster convergence

for very large data
each step is cheap & can be done completely in parallel

Convergence: monitor $\alpha$'s

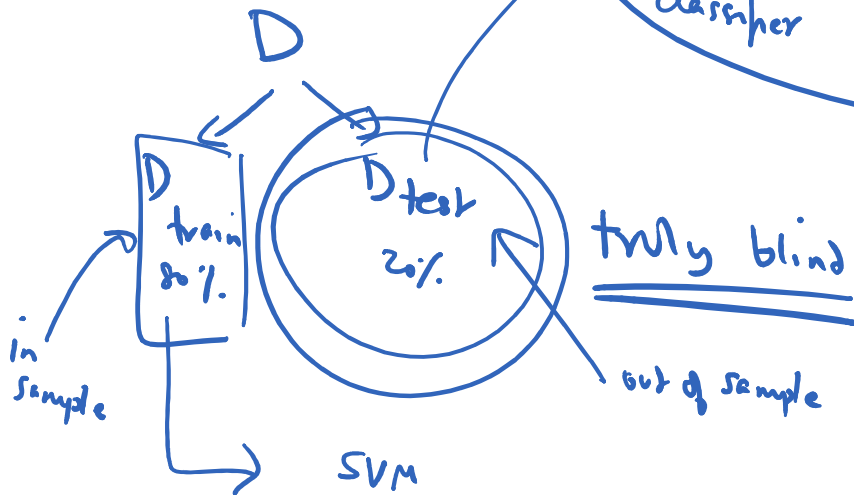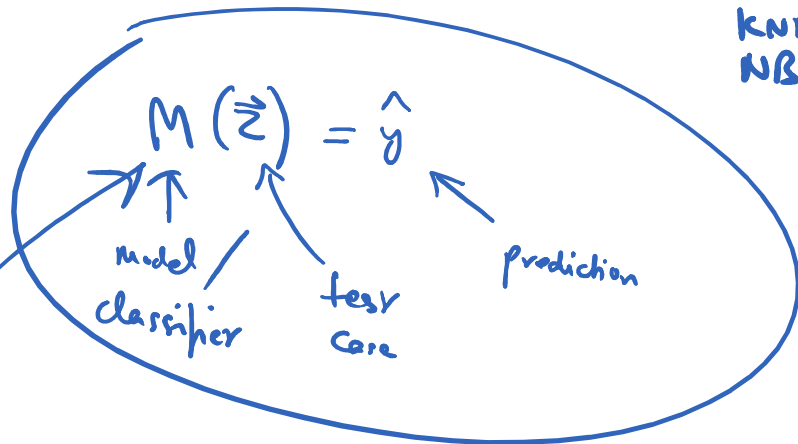$$\| \alpha_{prev} - \alpha_{new} \| \le \varepsilon \qquad Stop$$

↑
Convergence threshold

# Classification Evaluation

$D \leftarrow$ Sample of labeled data

$$\left( \vec{x}_i, y_i \right)^n_{i=1}$$

point — true class

SVM
DT
KNN
NB

$$M(\vec{z}) = \hat{y}$$

model / classifier — test case — prediction

$D$

$D_{train}$ 80%  — in sample

$D_{test}$ 20% — truly blind — out of sample

$\vec{z}$, true class $y$

$\hat{y}$ vs $y$

SVM
$\rightarrow C?$ $\leftarrow$ how to chose C

KNN
$\rightarrow k \leftarrow$ # of neighbors to select..

$D_{train}$ 90%

$D_{validation}$ 10% $\leftarrow$ used for parameter tuning

$$Accuracy = \frac{\# \text{ of correct predictions}}{|D_{test}|}$$

$\begin{array}{l} z_1 \rightarrow \hat{y}_1 \quad y_1 \\ z_2 \rightarrow \hat{y}_2 \quad y_2 \end{array}$

$$|D_{test}| = N \qquad = \frac{\#\{\hat{y}_i = y_i\}}{N}$$

$$|D_{test}|$$

$$
\begin{array}{c}
z_1 \rightarrow y_1 \quad y_1 \\
z_2 \rightarrow \hat{y}_2 \quad y_2 \\
\vdots \\
z_N \rightarrow \hat{y}_N \quad y_N
\end{array}
$$

predicted   true

Error Rate :   $1 - Accuracy$

$$: \quad \frac{\#\{\hat{y}_i \neq y_i\}}{N}$$

glolal measures of performance

---

$$
\begin{array}{cc}
\underline{P} & \underline{N} \\
10\% & 90\%
\end{array}
$$

$D_{train}$
$D_{test}$

$$\text{NULL}(\vec{z}) = N$$

NULL has 90% accuracy !!!

Claw specific accuracy / Coverage

$$
\begin{array}{cc}
\underline{P} & \underline{N} \\
0\% & 100\% \quad \longleftarrow \quad \text{accuracy}
\end{array}
$$

Confusion matrix
True

$$n_{pp} = \#\left( z \rightarrow \hat{y} = P \quad y = P \right)$$

|   | P | N |
|---|---|---|
| P | $n_{pp}$ | $n_{pn}$ |

$$\text{Prediction} \quad \begin{array}{c|c|c} & P & N \\ \hline P & n_{PP} & n_{PN} \\ \hline N & n_{NP} & n_{NN} \end{array}$$

TP: True positive

$$Z \to \hat{y}=N \quad y=N$$

TN: True negative



Predict $\quad \begin{array}{c|c|c} & P & N \\ \hline P & TP & \boxed{FP} \\ \hline N & \boxed{FN} & TN \end{array}$ $\quad$ True

FP: False Positive

$$Z \to \hat{y}=P \quad y=N$$

FN: False Negative

$$Z \to \hat{y}=N \quad y=P$$

Confusion matrix

$$\text{Predict} \quad \begin{array}{c|c|c|c} & P & N & \\ \hline P & \begin{array}{c}TP\\n_{PP}\end{array} & \begin{array}{c}FP\\n_{PN}\end{array} & m_P \\ \hline N & \begin{array}{c}FN\\n_{NP}\end{array} & \begin{array}{c}TN\\n_{NN}\end{array} & m_N \\ \hline & n_P & n_N & n \end{array}$$

$$|D_{test}| = n$$

$$n_P = \text{true \# of } P$$

$$n_N = \text{true \# of } N$$

$$m_P = \text{\# of predicted } P$$

$$m_N = \text{\# of predicted } N$$

## Positive class:

Precision or accuracy for P

$$= \frac{\text{\# of correct } P \text{ predictions}}{\text{\# of predictions for } P} = \frac{n_{PP}}{m_P} = \boxed{\frac{TP}{TP+FP}}$$

Coverage or recall for P

$$= \frac{\text{\# of correct } P}{\text{\# of } P \text{ in } D_{test}} = \frac{TP}{n_P} = \boxed{\frac{TP}{TP+FN}}$$

Negative class

Precision: $\quad \dfrac{TN}{-} = \dfrac{n_{NN}}{m_N}$

Precision : $\dfrac{TN}{TN + FN} = \dfrac{n_{NN}}{m_N}$

recall / coverage : $\dfrac{TN}{n_N} = \dfrac{TN}{TN + FP}$

## Precision vs. recall tradeoff

F score: harmonic mean of
P                 Precision & recall for a class

$$FScore_P : \dfrac{2}{\dfrac{1}{prec_P} + \dfrac{1}{recall_P}} = \boxed{\dfrac{2 \cdot prec_P \cdot recall_P}{prec_P + recall_P}}$$

Max = 1
min = 0