## PCA
principal
components
analysis

data dependent
dimensionality
reduction

- SVD    + applications
- least squares
- maybe – perceptron alg.

## Projects

— google form
— presenting et Microsoft?

PCA & SVD

have data set $x_1, \dots, x_m$ $x_i \in \mathbb{R}^n$

$$X = \begin{pmatrix} -\ x_1\ - \\ \vdots \\ -\ x_m\ - \end{pmatrix}$$

PCA: Fix $k$. find orthonormal vectors $\vec{v_1}, \dots, \vec{v_k}$

s.t. $x_i \simeq \sum_{j=1}^{k} (x_i, v_j) \vec{v_j}$

Find $v_1$ to min $\frac{1}{m} \sum_{i=1}^{M} \text{dist}\left(x_i, \text{line}(v_1)\right)^2$

$\equiv$ max $\frac{1}{m} \sum_{i=1}^{M} (x_i, v_1)^2$ variance.

$v_1$ principal eigenvector of matrix $X^T X$

$$X^T X = Q\ D\ Q^T \quad \longleftarrow \qquad Q Q^T = Q^T Q = I$$

↑
orthogonal.

$$Q = \begin{pmatrix} | & & | \\ q_1 & \cdots & q_n \\ | & & | \end{pmatrix}$$

symmetric matrix

$q_1, \dots, q_n$ eigenvectors
$\lambda_1, \dots, \lambda_n$

$$D = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$$

pos semi-definite.
all eigenvalues nonnegate

$$X^T X \,(q_1) = \lambda_1 q_1$$

$q_1$ best choice for $v_i$

singular value decompostn   SVD

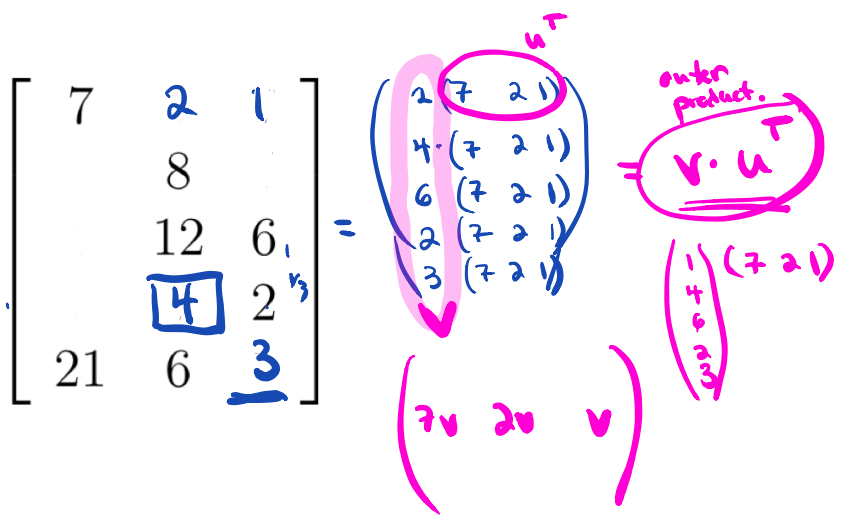$\Rightarrow$ tells us best way to approximate
our matrix with a "low rank"
matrix

movies

$$\text{people} \begin{bmatrix} 7 & ? & ? \\ ? & 8 & ? \\ ? & 12 & 6 \\ ? & ? & 2 \\ 21 & 6 & ? \end{bmatrix}$$

can we reconstruct missing entries?

Suppose ~~assume~~ that this matrix is rank 1

every row is a multiple of every other row.

$$\begin{bmatrix} 7 & 2 & 1 \\ & 8 & \\ & 12 & 6 \\ & \boxed{4} & 2 \\ 21 & 6 & \underline{3} \end{bmatrix} = \begin{pmatrix} 2 & (7 & 2 & 1) \\ 4 \cdot (7 & 2 & 1) \\ 6 & (7 & 2 & 1) \\ 2 & (7 & 2 & 1) \\ 3 & (7 & 2 & 1) \end{pmatrix}$$

$u^T$

outer product.

$$v \cdot u^T$$

$$\begin{pmatrix} 1 \\ 4 \\ 6 \\ 2 \\ 3 \end{pmatrix} (7 \ 2 \ 1)$$

$$\begin{pmatrix} 7v & 2v & v \end{pmatrix}$$

**Rank 0** matrix    all 0.

**Rank 1**

$$\mathbf{A} = \mathbf{u}\mathbf{v}^\top = \begin{bmatrix} - & u_1\mathbf{v}^\top & - \\ - & u_2\mathbf{v}^\top & - \\ & \vdots & \\ - & u_m\mathbf{v}^\top & - \end{bmatrix} = \begin{bmatrix} | & | & & | \\ v_1\mathbf{u} & v_2\mathbf{u} & \cdots & v_n\mathbf{u} \\ | & | & & | \end{bmatrix}$$

**Rank 2.**

$$\mathbf{A} = \mathbf{u}\mathbf{v}^\top + \mathbf{w}\mathbf{z}^\top = \begin{bmatrix} - & u_1\mathbf{v}^\top + w_1\mathbf{z}^\top & - \\ - & u_2\mathbf{v}^\top + w_2\mathbf{z}^\top & - \\ & \vdots & \\ - & u_m\mathbf{v}^\top + w_m\mathbf{z}^\top & - \end{bmatrix} = \begin{bmatrix} | & | \\ \mathbf{u} & \mathbf{w} \\ | & | \end{bmatrix} \cdot \begin{bmatrix} - & \mathbf{v}^\top & - \\ - & \mathbf{z}^\top & - \end{bmatrix}$$
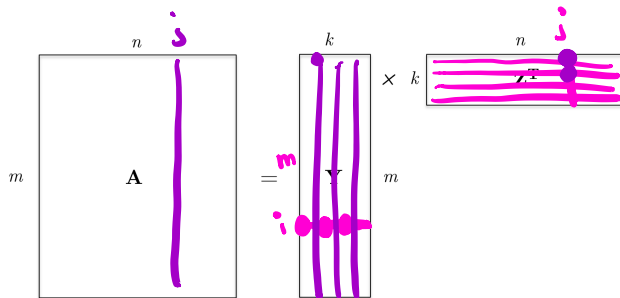


Figure 1: Any matrix $\mathbf{A}$ of rank $k$ can be decomposed into a long and skinny matrix times a short and long one.
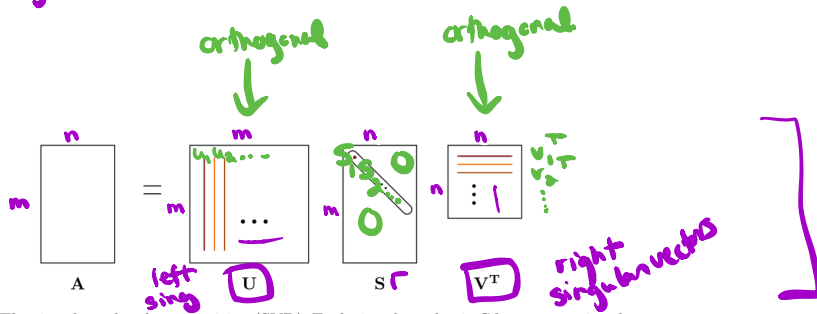
SVD of a matrix $m \times n$ matrix

orthogonal          orthogonal



n                    m                    n                    n  $V^T$
m  [A]  =  m  [U]          m  [S]     n  [$V^T$]
                left                    S           right
                sing                                singular vectors

Figure 2: The singular value decomposition (SVD). Each singular value in **S** has an associated left singular vector in **U**, and right singular vector in **V**.

diag entries
of
**S**
are called
singular values

$$A = \sum_{i=1}^{\min(m,n)} s_i u_i v_i^T$$

$$S_1 \geq S_2 \geq \dots \geq 0$$

can be computed in $\min\left(O(m^2 n), O(n^2 m)\right)$ time.

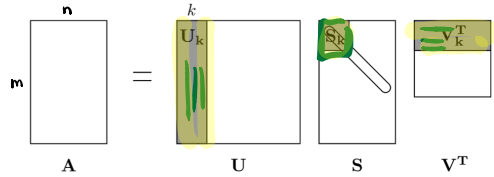Suppose we want the "best" rank $k$
approx to A

Figure 3: Low rank approximation via SVD. Recall that **S** is non-zero only on its diagonal, and the diagonal entries of $S$ are sorted from high to low. Our low rank approximation is $\mathbf{A}_k = \mathbf{U}_k\mathbf{S}_k\mathbf{V}_k^\top$.

$$A_k = m \left[\, U_k \,\right] \cdot \left[\, S_k \,\right] \cdot \left[\, V_k^T \,\right]$$

**Thm**

This low-rank approx is optimal in the sense that $\forall$ matrix $A$ $(m \times n)$ and rank target $k \geq 1$ and any other rank $k$ matrix $B$ $(m \times n)$

$$\|A - A_k\|_F^2 \leq \|A - B\|_F^2$$

$$\|M\|_F^2 = \sum_i \sum_j m_{ij}^2$$

$$\sum_{ij}(a_{ij} - a_{ij}^k)^2 \leq \sum_{ij}(a_{ij} - b_{ij})^2$$

Frobenius

Relationship between PCA & SVD.    $(AB)^T = B^T A^T$

$$X^T X = Q D Q^T$$

$$X = U S V^T$$

$q_1 = v_1$
$\lambda_1 = S_1^2$

$$X^T X = (USV^T)^T USV^T$$
$$= V^T S^T U^T U S V^T$$
$$\underbrace{U^T U}_{I}$$

$$X^T X = V S^2 V^T$$

$$\boxed{XX^T = U S^2 U^T}$$

$XX^T_{ij}$

$X = \text{custom} \left( \left\| \; \right\|_\text{products} \right)$

$$x_i = \sum_{j=1}^{k} (x_i, v_j) \vec{v_j}$$
$i^{th} \text{ row}$

---

## Application 1:    Denoising.

Suppose $A$ is a rank $k$ matrix.

$$C = A + \boxed{N}$$
$\quad\quad\quad\quad\quad$ noise matrix $\quad$ each entry of $N$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ is indep $\; N(0, \sigma^2)$

Then claim is    if variance of noise
Sufficiently small.

then    $\| C_k - A \|_F^2$   small w.h.p.
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ ↑
$\quad\quad\quad\quad\quad\quad\quad\quad$ depend on variance

$$A = \sum_{i=1}^{k} S_i u_i v_i^T \quad\quad \text{in } C \quad\quad\quad S_j \ll S_1 \dots S_k$$
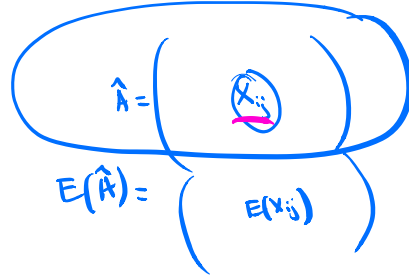$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad j > k$

**Thm:** $\hat{A}$ $m \times n$ matrix of indep r.v.'s whose variances are bounded by $\sigma^2$

If $\boxed{A = E(\hat{A})}$
is rank $k$

then w.h.p.

$$\left\| A - \hat{A}_k \right\|_F^2 = O\left( k \sigma^2 (m+n) \right)$$

$\hat{A} = \begin{pmatrix} X_{ij} \end{pmatrix}$

$E(\hat{A}) = \begin{pmatrix} E(X_{ij}) \end{pmatrix}$

$X_{ij} = \begin{array}{cc} 1 & \frac{2}{3} \\ 0 & \frac{1}{3} \end{array}$

$\dfrac{k\sigma^2(m+n)}{m \cdot n} = O(1)$

avg entry

$$\hat{A} = \underbrace{E(\hat{A})}_{\substack{A \\ \text{rank } k}} + \underbrace{(\hat{A} - E(\hat{A}))}_{\substack{\text{deviation from exp.} \\ \text{mean } 0 \\ \text{noise}}}$$

---

Collaborative Filtering
recommendations

movies

$\text{people} \begin{bmatrix} \boxed{3} \;\boxed{2}\; \boxed{1}\; 3 \\ 2 \; 4 \; 4 \; 5 \\ 1 \; 3 \; 1 \; 4 \end{bmatrix} \longrightarrow \begin{bmatrix} 3 \; 2 \; ? \; 3 \\ 2 \; ? \; 4 \; ? \\ ? \; 3 \; ? \; 4 \end{bmatrix}$

$p_{ij} = p$

$R$ $\qquad$ $R^?$

ground truth.

assumption: $R$ is rank $k$

human victory

Honey I Shrunk movie

$R$

$n$ movies



people $m$

humor
violence
romance
indie
⋮

$R$ is rank $k$
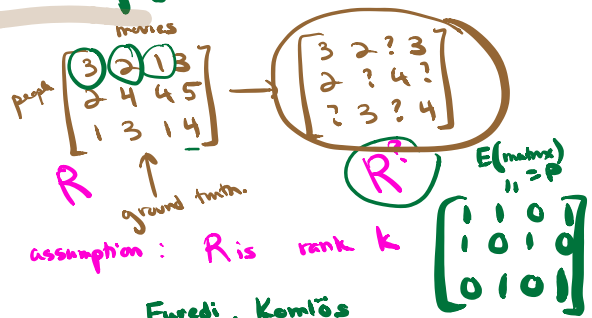Assume that there is $P = \begin{bmatrix} p_{ij} \end{bmatrix}$

s.t. prob that there is a rating available for that entry is $p_{ij}$

Let's assume we know $P$.

movies
$\begin{bmatrix} 3 & 2 & 1 & 3 \\ 2 & 4 & 4 & 5 \\ 1 & 3 & 1 & 4 \end{bmatrix}$ people $\longrightarrow$ $\begin{bmatrix} 3 & 2 & ? & 3 \\ 2 & ? & 4 & ? \\ ? & 3 & ? & 4 \end{bmatrix}$

$R$
ground truth.

$\hat{R}$?

$R$?

assumption: $R$ is rank $k$

$E(\text{matrix})_{ij} = P$

$\begin{bmatrix} 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}$

Define
$$\hat{R} = \begin{cases} \dfrac{R_{ij}}{P_{ij}} & \text{if entry } (i,j) \text{ is present} \\ 0 & \text{o.w.} \end{cases}$$

$E(\hat{R}_{ij}) = R_{ij}$
$= P_{ij} \dfrac{R_{ij}}{P_{ij}} + (1-P_{ij})0$
$= R_{ij}$

$\Longrightarrow$ $\hat{R}_k$ is very close to $R$.

Furedi, Komlós

Thm: $\hat{A}$ $m \times n$ matrix of indep r.v.'s whose variances are bounded by $\sigma^2$

If $A = E(\hat{A})$ is rank $k$
then w.h.p.

$\| A - \hat{A}_k \|_F^2 = O(k\sigma^2 (m+n))$

$\hat{A} = A + \underbrace{(\hat{A} - A)}_{\text{mean } 0 \text{ variance.}}$

$\hat{A} = \begin{pmatrix} \hat{x}_{ij} \end{pmatrix}$
$E(\hat{A}) = \begin{pmatrix} E(x_{ij}) \end{pmatrix}$

$x_{ij} = \begin{cases} 1 & \frac{2}{3} \\ 0 & \frac{1}{3} \end{cases}$

Assume $P$ itself is low rank

given $R$?
construct
matrix $\hat{P}$

whose $(i,j)$ entry is 1
see rating there and 0
o.w.

movies

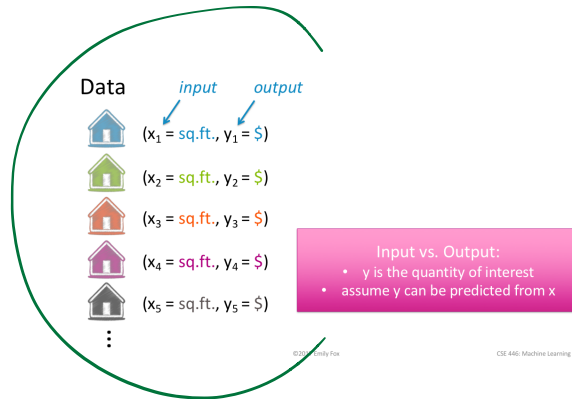people $\left(\ \ \right)\left(\ \ \right)$

$E(\hat{P}) = P$

$\Rightarrow \widehat{\hat{P}}$ very close to $P$

*Linear regression.*

How much is my house worth?

I want to list my house for sale

3  ©2017 Emily Fox

$f(\text{square footage}) \rightarrow \text{price}$

Model –
How we *assume* the world works

$\varepsilon_i: f(x)$

y
price ($)

square feet (sq.ft.)  X

6  ©2017 Emily Fox

Model –
How we *assume* the world works

y
price ($)

"Essentially, all models are wrong, but some are useful."
George Box, 1987.

square feet (sq.ft.)  X

7  ©2017 Emily Fox    CSE 446: Machine Learning

Simple linear regression model

y
price ($)

$y_i = w_0 + w_1 x_i + \varepsilon_i$

$f(x) = w_0 + w_1 x$

square feet (sq.ft.)  X

8  ©2017 Emily Fox

$(x_i, y_i)$ pairs.

Find best choice for $w_0, w_1$

## Add more inputs

$f(x) = w_0 + w_1$ sq.ft. $+ w_2$ #bath

y

price ($)

# bathrooms

x[2]

square feet (sq.ft.)  x[1]

## Many possible inputs

- Square feet
- # bathrooms
- # bedrooms
- Lot size
- Year built
- ...

## "Cost" of using a given line

y

Residual sum of squares (RSS)

price ($)

$w_0, w_1$

$$RSS(w_0, w_1) = \sum_{i=1}^{N} (y_i - [w_0 + w_1 x_i])^2$$

square feet (sq.ft.)   $x_i$   x

Given $(x_i, y_i)$
$1 \leq i \leq N$

Problem: find $w_0 \& w_1$
to minimize

# RSS for multiple regression

Linear Regression

$(x_i, y_i)$ pairs.

PCA-

y
price ($)

square feet (sq.ft.)    x

y
price ($)

square feet (sq.ft.)    x

$$m \left( \quad \right) k \left( \quad \right)$$

$k$ $n$

$k(m+n)$

$m \cdot n$

Supervised learning. & Perceptron Algorithm

labelled data:           use that to come up
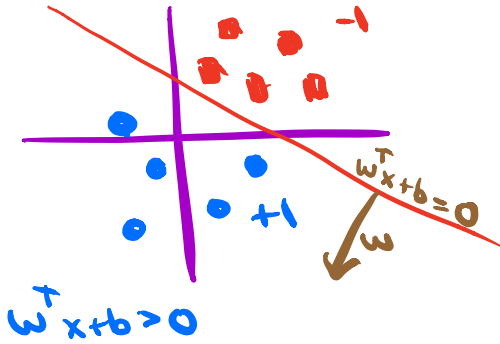                    with a   way to  give answers  on
                    data    haven't  seen.  $(\vec{x_i}, y_i)$

$w^T x + b < 0$



$w^T x + b = 0$

$w^T x + b > 0$

binary classification.
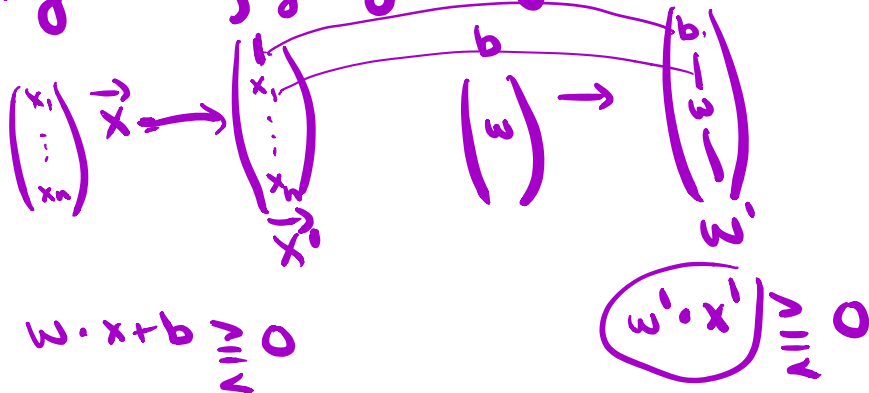label for each pt $\{-1, 1\}$

Assumption:
data set $\underline{\underline{is}}$   linearly separable.

$\exists \; \vec{w} \; \& \; b \quad s.t.$

$y_i = 1 \quad \Rightarrow \quad w^T x + b > 0$

$y_i = -1 \quad \Rightarrow \quad w^T x + b < 0$

Objective:  find one

Simplify life  by getting rid of   additive const b.

$\begin{pmatrix} x_i \\ \vdots \\ x_n \end{pmatrix} \vec{x} \rightarrow \begin{pmatrix} 1 \\ x_i \\ \vdots \\ x_n \end{pmatrix} \vec{x}' \qquad \begin{pmatrix} b \\ w \end{pmatrix} \rightarrow \begin{pmatrix} b \\ | \\ w \\ | \end{pmatrix} w'$

$w \cdot x + b \gtreqless 0$          $\boxed{w' \cdot x' \gtreqless 0}$

$y = 1 \qquad\qquad w^T x > 0$
$y = -1 \qquad\qquad w^T x < 0$

make a mistake if   $\boxed{y_i \, w^T x_i \leq 0}$

```
Initialize w⃗ = 0⃗                          // Initialize w⃗. w⃗ = 0⃗ misclassifies everything.
while TRUE do                             // Keep looping
    m = 0                                  // Count the number of misclassifications, m
→ for (xᵢ, yᵢ) ∈ D do                      // Loop over each (data, label) pair in the dataset, D
      if yᵢ(w⃗ᵀ · x⃗ᵢ) ≤ 0 then    made mistake    // If the pair (x⃗ᵢ, yᵢ) is misclassified
        w⃗ ← w⃗ + yᵢx⃗ᵢ  ←                   // Update the weight vector w⃗
        m ← m + 1                          // Counter the number of misclassification
      end if
    end for
    if m = 0 then                          // If the most recent w⃗ gave 0 misclassifications
      break                                // Break out of the while-loop
    end if
end while                                  // Otherwise, keep looping!
```
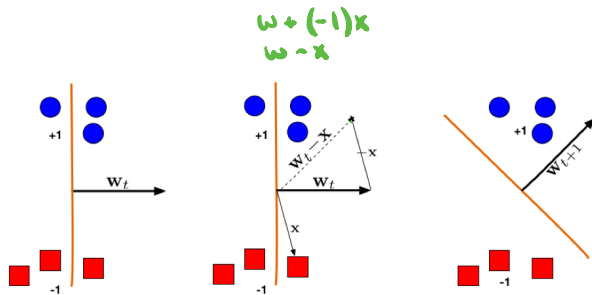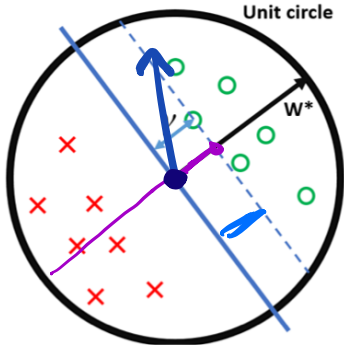
$$w + (-1)x$$
$$w - x$$



*Illustration of a Perceptron update. (Left:) The hyperplane defined by* $\mathbf{w}_t$ *misclassifies one red (-1) and one blue (+1) point. (Middle:) The red point* $\mathbf{x}$ *is chosen and used for an update. Because its label is -1 we need to **subtract** $\mathbf{x}$ from $\mathbf{w}_t$. (Right:) The udpated hyperplane $\mathbf{w}_{t+1} = \mathbf{w}_t - \mathbf{x}$ separates the two classes and the Perceptron algorithm has converged.*

$w^*$ to a correct answer

$$\exists \, w^* \qquad y_i \left( \text{sign} \left( x_i^T w^* \right) \right) > 0$$

wlog $\|w^*\| = 1$

to simplify

$$x_i' = \frac{x_i}{\max_j \|x_j\|}$$

$$\|x_i'\| \leq 1$$

**Unit circle**

$\gamma$ is called the margin
$$= \min_{(x_i, y)} |x_i^T w^*|$$

Thm: $m \le \dfrac{1}{\gamma^2}$

Look at 2 quantities.

$w^T w^*$ ① — increasing

$w^T w$ ② — can't be going up fast

---

When we make a mistake:
$y(w^{*T}x) \le 0$
$y((w^T x) > 0$

① $(w + yx)^T w^*$
$= w^T w^* + \underbrace{y \cdot x^T w^*}_{\ge \gamma}$

every mistake $\Rightarrow$ increase $w^T w^*$ by at least $\gamma$

② $(w + yx)^T (w + yx)$
$= w^T w + \underbrace{2y x^T w}_{\le 0} + \underbrace{y^2 x^T x}_{\le 1}$

every mistake, increase $w^T w$ by at most $1$

$0 < m\gamma \le \boxed{\dfrac{w^T w^*}{w^T w^*}}$

$= w^T w^*$

$\le \|w\| \|w^*\| \cos\theta$

$\le \|w\| \|w^*\|$
$\qquad\qquad \uparrow$
$\qquad\qquad 1$

$= \|w\|$

$= \sqrt{w^T \cdot w}$

$\le \sqrt{m}$

$m\gamma \le \sqrt{m}$

$m^2 \gamma^2 \le m$

$m \le \dfrac{1}{\gamma^2}$

Bad example for
perceptron alg