

Logistic regression

A probabilistic classification model

$$\{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^d \quad y_i \in \{0, 1\}$$

We just saw generative models

$$\text{e.g. } P(x|y=c, \theta) = N(x|\mu_c, \Sigma_c)$$

We then compute posterior $P(y|x, \theta)$

$$\propto P(y) P(x|y, \theta)$$

to classify.

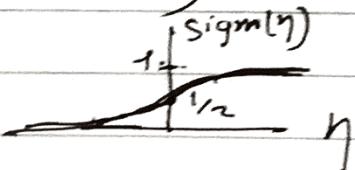
One might want to have discriminative models directly giving us the probability $P(y|x, \theta)$

We saw that two class Gaussian model with tied covariance matrices had a posterior class probability given by a sigmoidal / logistic function of a linear combination of ~~variables~~ variable. We could try to ~~learn~~ learn such a model directly from labeled data.

$$P(y|x, \theta) = \text{Ber}(y | \text{sigm}(\theta^T x))$$

remember

$$\text{sigm}(\eta) = \frac{1}{1 + e^{-\eta}}$$



BTW, if we need to add a constant shift, say $P(y|x, \theta) = \frac{1}{1 + e^{-(w^T x + w_0)}}$

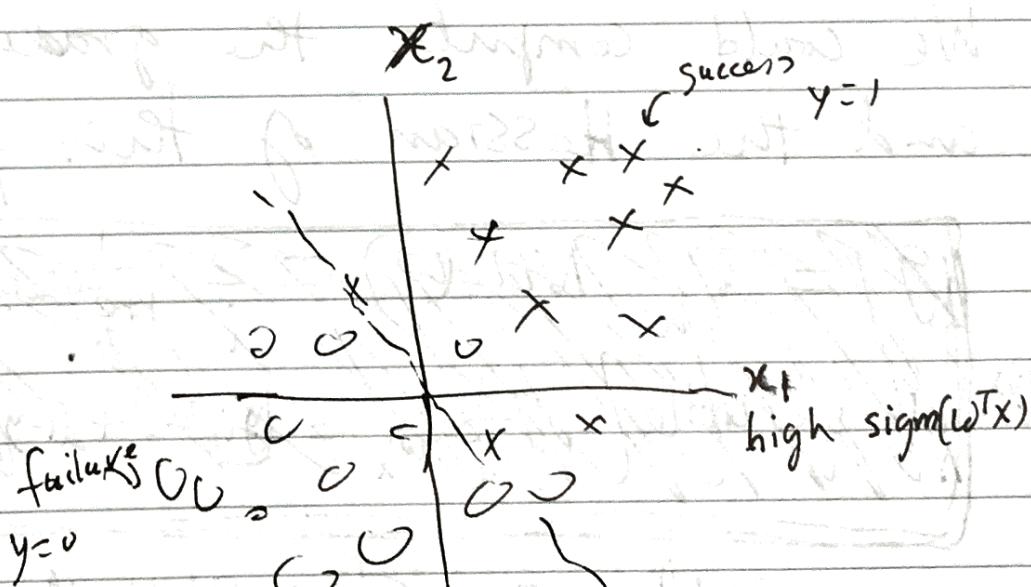
We could just define bigger vector

$$w \rightarrow (w_0, w) = \tilde{w}$$

$$x \rightarrow (1, x) = \tilde{x} \quad P(y|x, \theta) = \frac{1}{1 + e^{-\tilde{w}^T \tilde{x}}}$$

$$\prod_i P(y_i | x_i, w) = \prod_i \left[\mu_i^{I(y_i=1)} (1-\mu_i)^{I(y_i=0)} \right]$$

where $\mu_i = \frac{1}{1 + e^{-w^T x_i}}$



$$\text{high sigm}(w^T x)$$

MLE Negative Log likelihood

$$NLL(\omega) = - \sum_i [y_i \log \mu_i + (1-y_i) \log(1-\mu_i)]$$

$$\text{By the way } 1-\mu_i = \left(1 - \frac{1}{1+e^{-\omega^T x}}\right) = \frac{e^{-\omega^T x}}{1+e^{-\omega^T x}}$$

$$\text{So, if we use } \tilde{y}_i = \pm 1 = y_{i-1} = \frac{1}{e^{\omega^T x} + 1}$$

$$\begin{aligned} NLL(\omega) &= - \sum_i \log \frac{1}{1+e^{\tilde{y}_i \omega^T x_i}} \\ &= - \sum_i \log \left(1 + e^{-\tilde{y}_i \omega^T x_i}\right) \end{aligned}$$

Unfortunately we can't wait down $\hat{\omega}$ in one shot
 We could compute the gradient
 and the Hessian of this function

$$\boxed{\begin{aligned} \nabla_w \mu_i &= \frac{\partial \mu_i}{\partial w} = \frac{\partial}{\partial w} \frac{1}{1+e^{-\omega^T x_i}} = \frac{e^{-\omega^T x_i}}{(1+e^{-\omega^T x_i})^2} \cdot (-\omega^T x_i) \\ \Rightarrow \nabla_w \mu_i &= \frac{-\omega^T x_i}{(1+e^{-\omega^T x_i})^2} \end{aligned}}$$

$$\begin{aligned} g &= \nabla_w NLL = - \sum_i [y_i \nabla_w \log \mu_i + (1-y_i) \nabla_w \log(1-\mu_i)] \\ &= - \sum_i [(1-\mu_i)y_i x_i + (1-y_i)\mu_i(-x_i)] \end{aligned}$$

$$f = \sum_i (\mu_i - y_i) x_i = X^T (\mu - y)$$

$$\begin{aligned} f &= \nabla_{\mu} g(\mu)^T = \sum_i (\nabla_{\mu} \mu_i) x_i^T \\ &= \sum_i \mu_i (1 - \mu_i) x_i x_i^T \\ &= X^T S X \end{aligned}$$

where $S \triangleq \text{diag}(\mu_1(1-\mu_1), \dots, \mu_N(1-\mu_N))$

$$X = \begin{pmatrix} x_{11} & \cdots & x_{1D} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{ND} \end{pmatrix}$$

$$X^T = \begin{pmatrix} x_{11} & \cdots & x_{N1} \\ \vdots & \ddots & \vdots \\ x_{1D} & \cdots & x_{ND} \end{pmatrix}$$

$$\mu - y = \begin{pmatrix} \mu_1 - y_1 \\ \vdots \\ \mu_N - y_N \end{pmatrix}$$

$$S = \begin{pmatrix} \mu_1(1-\mu_1) & & & \\ & \mu_2(1-\mu_2) & & \\ & & \ddots & \\ & & & \mu_N(1-\mu_N) \end{pmatrix}$$

How do we minimize $NLL(w)$?

One thought: ~~Steepest descent~~

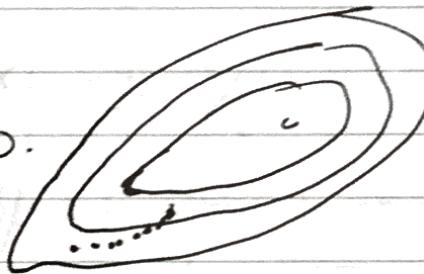
$$\theta_{k+1} = \theta_k - \eta_k g_k$$

Issues with step size

Big η
Potential oscillation



Small η
May be too slow.



There are many methods to deal with issues of stepsize selection. We describe ^{one} the popular methods in the context of ~~the~~ logistic regression.

Newton's method

Initialize θ_0 :
for $k = 1, 2, \dots$ until convergence do

Evaluate $g_k = \nabla f(\theta_k)$

Evaluate $H_k = \nabla^2 f(\theta_k)$;

Solve $H_k d_k = -g_k$;

Use line search to find step size γ_k along d_k ;

$$\theta_{k+1} = \theta_k + \gamma_k d_k$$

Roughly $f(\theta+d) = f(\theta) + g \cdot d$
 $+ \frac{1}{2} d^T H d$
+ ...

The second order approximation
is minimized when $g + H d = 0$
(assuming H_k is positive definite)

To avoid overstepping one uses
line search in the direction of d .

For one dimension

For logistic regression with $\gamma_k = 1$

$$w_{k+1} = w_k - H_k^{-1} g_k$$

$$= w_k + (X^T S_k X)^{-1} X^T (y - \mu_k)$$

$$= (X^T S_k X)^{-1} \left[(X^T S_k X) w_k + X^T (y - \mu_k) \right]$$

$$= (X^T S_k X)^{-1} X^T S_k [X w_k + S_k^{-1} (y - \mu_k)] = (X^T S_k X)^{-1} X^T Z_k$$

where $z_k \triangleq X w_k + S_k^{-1} (y - \mu_k)$
 which is called the working response

Note that w_{k+1} optimizes

$$\sum_{ki} S_{ki} (z_{ki} - w^T x_i)^2$$

$$\nabla_w \sum S_i^2 = 0 \rightarrow \sum_{ki} [x_i^T S_{ki} z_{ki} - x_i^T S_{ki} x_i w] = 0$$

$$z_{ki} = w_k^T x_i + \frac{y_i - \mu_{ki}}{\mu_{ki}(1-\mu_{ki})}$$

Iterating recomputation of z_{ki}
 with least square minimization
 gives iteratively reweighted
 least square (IRLS) algorithm.

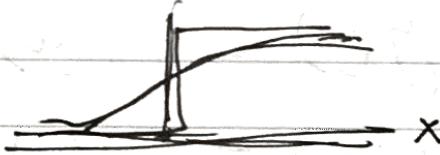
Now let us discuss issues of
 regularization. What if training
 data ~~was~~ for the two classes

were linearly separable, by chance?

$$\begin{matrix} 0 & 0 & 1 & x & x & + \\ 0 & 0 & 1 & x & & \\ 0 & 0 & 1 & x & x & \end{matrix}$$

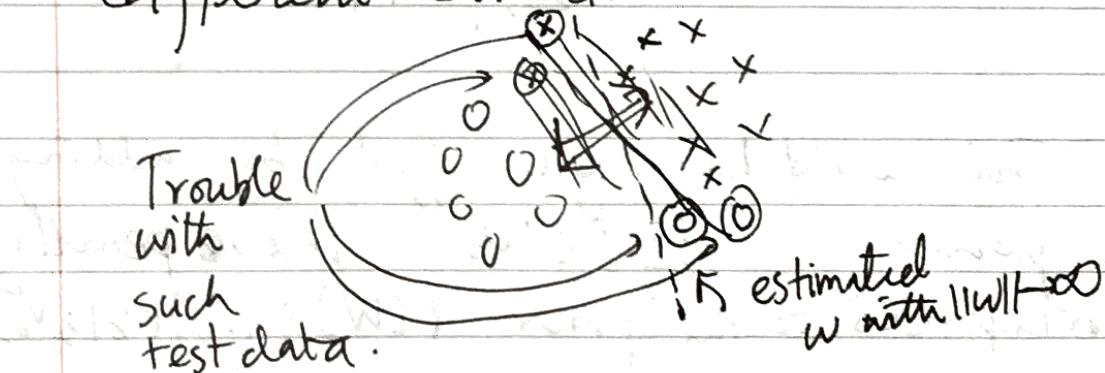
In that case,

once the direction
is established,
 $\|w\| \rightarrow \infty$



That way all $y=1$ have probability 1,
and all $y=0$ have probability 0,
maximizing the likelihood $\prod_{i:y=1} 1 \prod_{i:y=0} (1-0)$

This is brittle. Let's say the true model has finite w and a slightly different direction



One option is to add an L_2 regularization as in ridge regression

$$f'(\omega) = NLL(\omega) + \lambda \omega^T \omega$$

$$g'(\omega) = g(\omega) + 2\lambda \omega$$

$$H'(\omega) = H(\omega) + 2\lambda I$$

Also, one could generalize logistic regression to multiclass ~~setting~~ setting

$$P(y=c|x, \omega) = \frac{\exp(\omega_c^T x)}{\sum_{c'=1}^C \exp(\omega_{c'}^T x)}$$

This could be regularized by adding a penalty $\sum_c \omega_c^T v_0^{-1} \omega_c$. This penalty corresponds to the prior $P(\omega) = \prod_c N(\omega_c | 0, v_0)$

As we will see, just regularizing may not be enough!

Bayesian logistic regression

As we discussed before, instead of using the mode of $p(w|D)$, namely MAP estimate, we might want to integrate w over w .

let us take a more general look at Bayesian integration in a particular approximation that becomes accurate in the large data size limit.

Think of $P(\theta|D) = \frac{e^{-E(\theta)}}{Z}$ with
 $Z = \int d\theta' e^{-E(\theta')}$. If $E(\theta) = -\log P(\theta|D)$

$$Z = P(D) \triangleq \text{evidence}$$

Let us say θ^* is the MAP estimate (lowest energy \cong highest posterior prob.)

$$E(\theta) = E(\theta^*) + (\theta - \theta^*)^T g + \frac{1}{2} (\theta - \theta^*)^T H (\theta - \theta^*)$$

If θ^* is the mode, $g = 0$

$$\begin{matrix} g \in \mathbb{R}^D \\ H \in \mathbb{R}^{D \times D} \end{matrix}$$

$$Z = \int e^{-E(\theta)} d\theta = e^{-E(\theta^*)} \int d\theta e^{-\frac{1}{2}(\theta - \theta^*)^T H(\theta)}$$

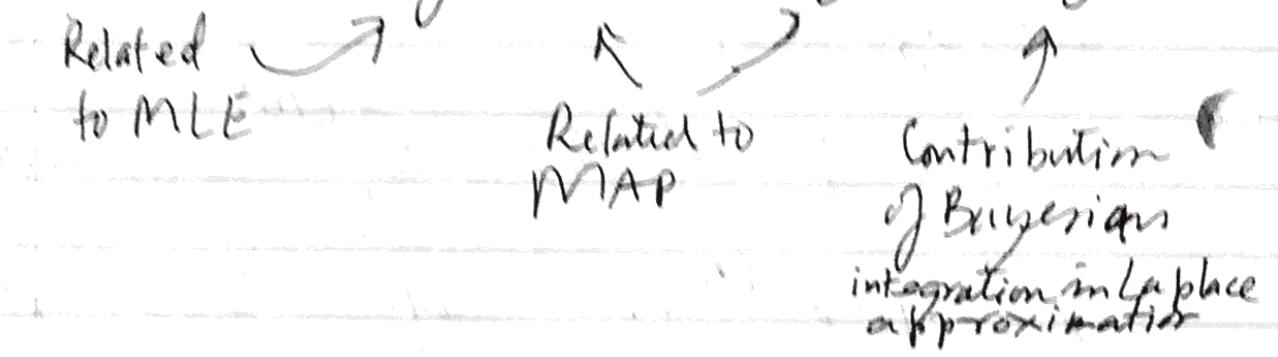
$$= e^{-E(\theta^*)} (2\pi)^{D/2} |H|^{-1/2}$$

This $\mathcal{N}(\theta^*)$ approximation is known as the Laplace approximation.

$$\log P(\theta) = -E(\theta^*) - \frac{1}{2} \log |H| + \text{constant}$$

$$= \boxed{\log P(\theta^*/\theta)} + \log P(\theta^*) - \frac{1}{2} \log |H| + \text{const.}$$

$$= \log P(\theta/\theta^*) + \log P(\theta^*) - \frac{1}{2} \log |H| + \text{const.}$$



Relation to Bayesian information criterion (BIC)

$$N \text{ pieces of data: } H = \sum_i H_i, H_i = \nabla D \log P(\theta, i)$$

If each $H_i \approx \bar{H}$, $H \approx N \bar{H}$

$$\log |H| \approx \log |N \bar{H}| = \log N^D \bar{H}^1$$

$$\log P(\theta) \sim \log P(\theta/\theta^*) + \log P(\theta^*) - \frac{D}{2} \log N \frac{1}{2} \log \bar{H}^1 + \dots$$

$\sim N$ ~ 1 $\sim \log N$ ~ 1

$$\sim \log P(D|\theta^*) - \frac{D}{2} \log N$$

For comparing models, the term $-\frac{D}{2} \log N$

[~~pro~~] penalizing more complex model -

(~~better~~ ^{higher} $\log P(D|\theta^*)$ achieved by bigger D offset by $-\frac{D}{2} \log N$)

Akaike information criterion (AIC)
considers $2D - 2\log P(D|\theta^*)$
 $= -2(\log P(D|\theta) - D)$

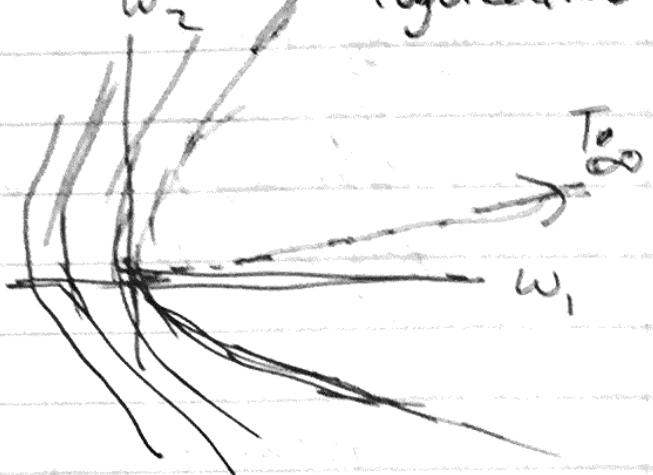
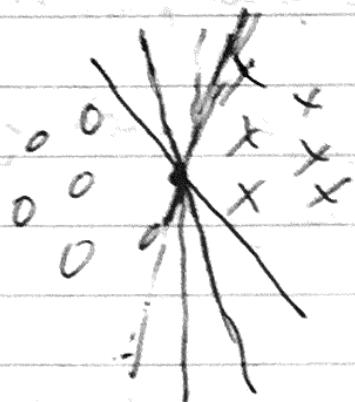
BIC depends on data size. Since $\frac{\log N}{2}$ is often greater than 1, BIC penalizes complex models more severely.

[Many prefer writing ~~KD~~ L
 $KD - 2\log p(D|\theta^*)$ and minimize it.]

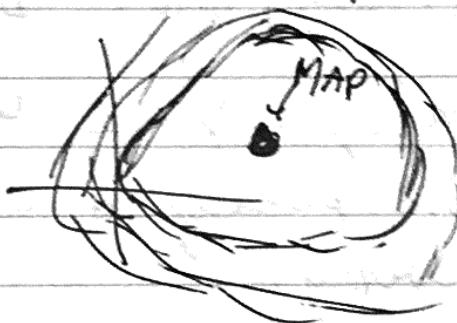
$K = \log N$ for BIC and $K=2$ for AIC]

Back to the example of perfectly separable case. If data comes really

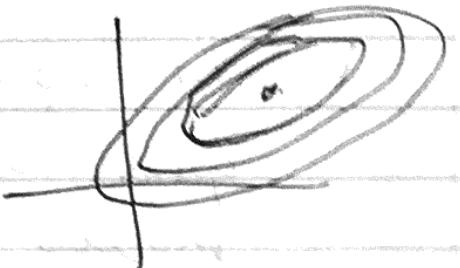
(as a result)
from p dist + logistic regression calling
the class (we will do a two class example),
we will see this phenomenon only when
 N is not too large w_2 loglikelihood



log Posterior
(with quadratic penalty)



Laplace approximation



Approximating posterior predictive

$$P(y|x, \bar{w}) = P(y|x, w) P(w|\bar{w}) \approx P(y|x, \hat{w})$$

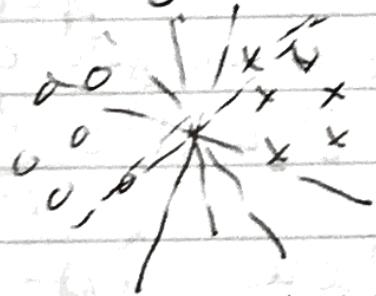
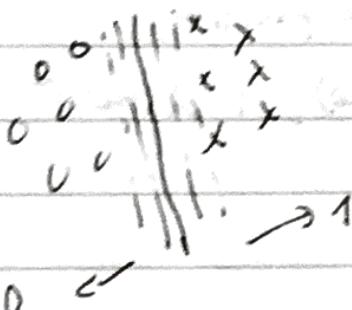
\hat{w}_{MAP} or $E[w]$

Alternative : Sample w from the posterior
Monte Carlo (MC)

$$P(y|x, \bar{w}) \approx \frac{1}{S} \sum_{s=1}^S \text{sigm}((\bar{w}^s)^T x)$$

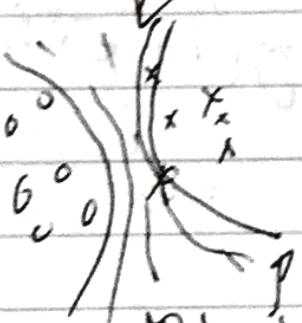
$$P(y_i|x_i, \hat{w})$$

Bayesian average over



The variability over orientation is important
either Laplace approx or M C

$$P(z_i|x_i, \hat{w})$$



regions of uncertainty

Online logistic regression and perceptrons

For optimizing $f(\theta) = \frac{1}{N} \sum_{i=1}^N f(\theta, z_i)$ $\rightarrow -\log P(y_i|x_i, \theta)$

one could take an online approach: change θ in response to new data.

Grad descent $\theta_k = \theta_{k-1} - \eta_k g_k$

for logistic reg: this means

$$w_k = w_{k-1} - \eta_k (y_i - g_i) x_i \quad y_i \in \{0, 1\}$$

If we consider using $y_i \in \{-1, 1\}$, and use $g_i = \hat{x}_i - 1 \approx \pm 1$
and $\eta_k = \frac{1}{2} \eta_k$ $w_k = w_{k-1} + \eta_k y_k x_i$ when $y_k \neq g_k$

Perceptron algorithm: $w_k = w_{k-1} + \eta_k y_k x_i$ otherwise

Support vector Machines (SVM's)

L_2 regularized empirical risk function

$$J(\omega, \lambda) = \sum_{i=1}^N L(y_i, \hat{y}_i; \omega) + \lambda \|\omega\|^2$$

where $\hat{y}_i; \omega = \omega^T x_i + \omega_0$

If $L(y, \hat{y}) = \sum_i (y_i - \hat{y}_i)^2 \Rightarrow$ ridge regression

" $L(y, \hat{y}) = \log(1 + e^{-y\hat{y}}) \Rightarrow$ logistic regression

$$\boxed{\text{redacted}} \quad y = +1, -1$$

We have seen that for ridge regression

$$\hat{\omega} = (X^T X + \lambda I_N)^{-1} X^T y \quad \boxed{\text{redacted}}$$

$$\hat{y} = \hat{\omega}^T x \quad \boxed{\text{redacted}} = Y^T X (X^T X + \lambda I_N)^{-1} X^T y = Y^T (X X^T + \lambda N I_N)^{-1} X^T y$$

This could be reexpressed as $\hat{y} = \sum_i \alpha_i K(x_i, x)$

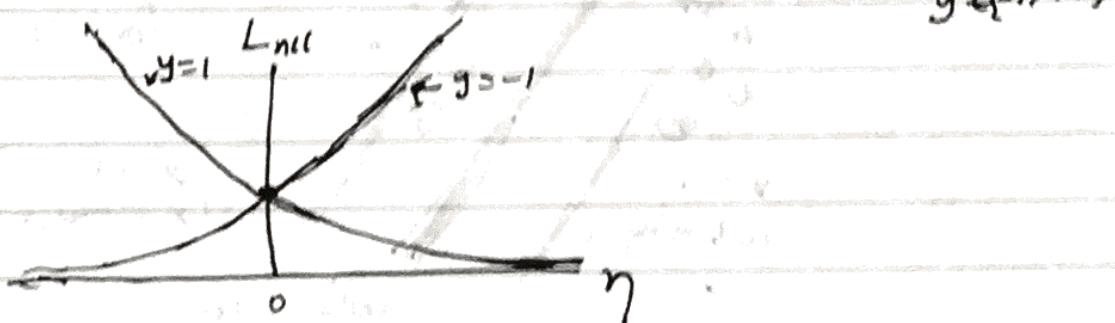
$$\alpha = Y \quad K(x_i, x_j) = \boxed{\text{redacted}} x_i^T x_j \quad \text{with } A = \boxed{\text{redacted}} X^T X + \lambda N I_N$$

Of course, here all x 's and y 's contribute

We will soon generalize the Kernel structure.

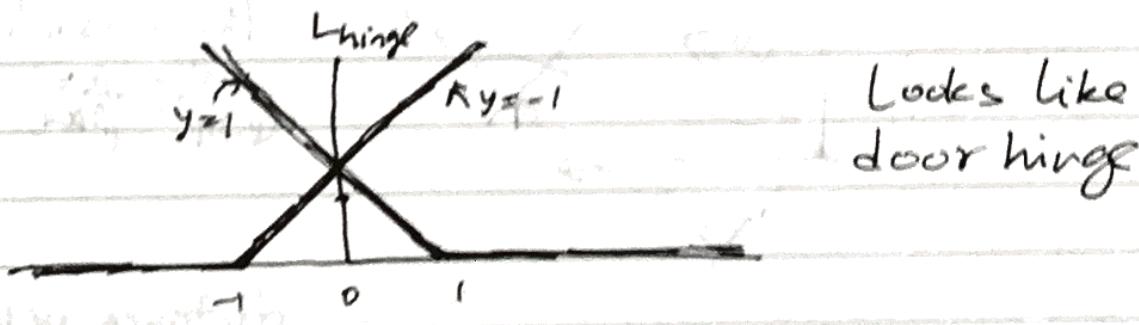
Coming off discussions of logistic regression
 let us first take up support vector classifier

Logistic regression: $\eta = f(x) = w_0 + \omega^T x$
 $L_{NLL}(y, \eta) = -\log p(y|x, \omega) = \log(1 + e^{-y\eta})$



Replace this loss function by hinge loss

$$L_{\text{hinge}}(y, \eta) = \max(0, 1 - y\eta) = (1 - y\eta)_+$$



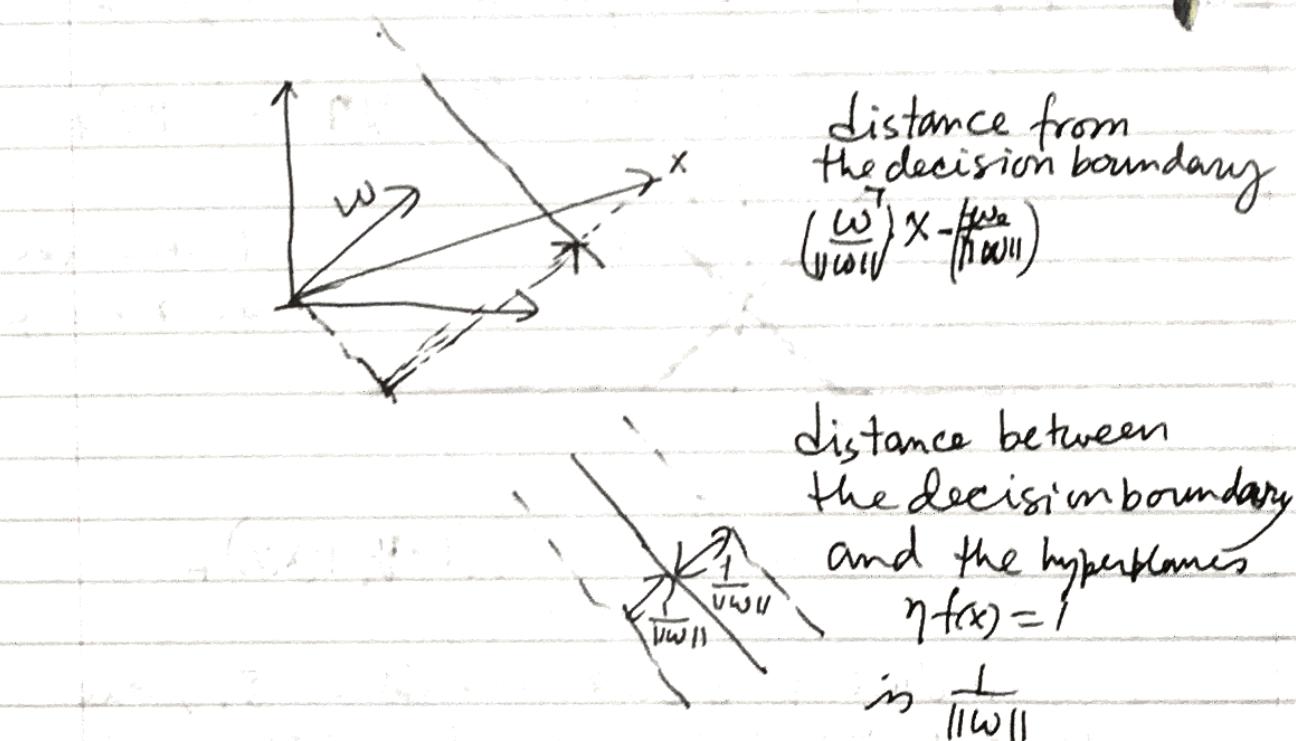
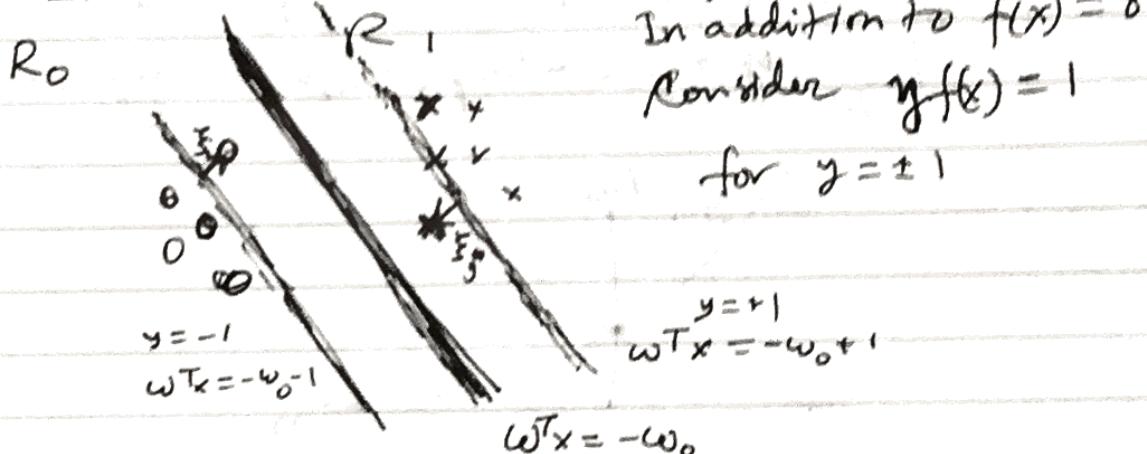
$$\min_{w, w} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N (1 - y_i f(x_i))_+$$

Introducing slack variables ξ_i , one could rewrite this optimization problem as

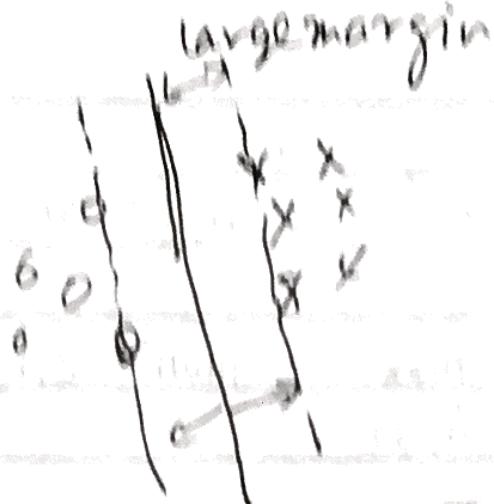
$$\min_{w, w, \xi} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \quad \text{st. } \xi_i \geq 0 \text{ and } y_i(w + w^T x_i) \geq 1 - \xi_i \text{ for } i = 1, \dots, N$$

The points are classified by $f(x) \geq 0$.
 $f(x) = w_0 + w^T x$ is the discriminant function.

Let consider the geometry of relevant hyperplane



If the ~~classes~~ classes are linearly separable and we have a large C we want all $\eta = 0$



Not



We want $\|w\|$ to be
as large s.t.

$$y_i(w_0 + w^T x_i) \geq 1 \text{ for all } i$$

That is equivalent to

$$\min \frac{1}{2} \|w\|^2, \text{ s.t. } y_i(w_0 + w^T x_i) \geq 1 \text{ for } i=1, \dots, N$$

In cases where perfect separation is not possible, we have the soft margin version

$$\min_{w, w_0, \xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i \quad \text{s.t. } y_i(w_0 + w^T x_i) \geq 1 - \xi_i$$

For the sake of simplicity let us discuss perfectly separable case.

First ~~w~~ $\min_w \|w\|^2 \quad \text{s.t. } y_i(w_0 + w^T x_i) \geq 1$

Could be rewritten using Lagrange multipliers.

Consider

$$L(w_0, w, \lambda) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \lambda_i (y_i (w_0 + w^T x_i) - 1)$$

Turns out $\min_{w, w} \max_{\lambda; \geq 0} \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \lambda_i (y_i (w_0 + w^T x_i) - 1)$

The key difference is the lagrange multipliers are sign restricted.

Think of it this way: If any constraint is not satisfied, ie. $y_i (w_0 + w^T x_i) - 1 < 0$

We can send $\lambda_i \rightarrow +\infty$ and

the obj. func. L would by ∞ . When minimizing over w , we will reject such solutions.

If a constraint is strictly satisfied

$$y_i (w_0 + w^T x_i) - 1 > 0$$

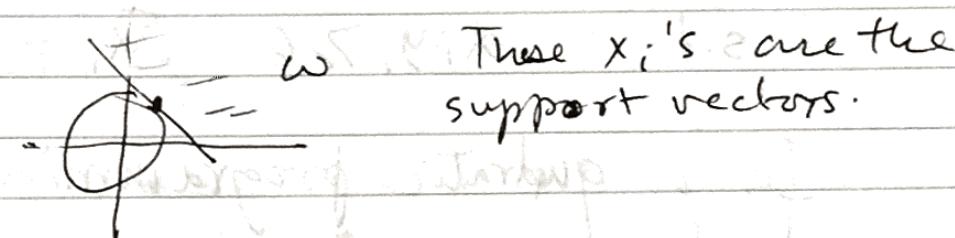
max over λ_i , sends $\lambda_i = 0$

So inactive constraints have corresponding $\lambda_i = 0$

What if $y_i (w_0 + w^T x_i) - 1 = 0$?

The margin case is interesting

Now, λ_i is determined by w variation being zero condition. It would usually be nonzero.



These x_i 's are the support vectors.

Turns out that in this case, we can switch the order of min & max

$$\max_{\lambda_i \geq 0} \min_{w_0 + w} \frac{1}{2} \|w\|^2 - \sum_i y_i ((w_0 + w \cdot x_i) - 1)$$

Since $w_0 + w$ are unrestricted, we could just differentiate.

$$w = \sum_i \lambda_i y_i x_i = \sum_i \overbrace{\lambda_i y_i}^{\alpha_i} x_i$$

$$0 = \sum_i \lambda_i y_i = \sum_i \overbrace{\lambda_i y_i}^{\alpha_i}$$

$$\tilde{L}(\alpha) = -\frac{1}{2} \sum_i \alpha_i x_i^\top x_i + \sum_i \alpha_i y_i$$

$\alpha_i y_i \geq 0$ for all i

and $\sum_i \alpha_i = 0$

$$\min_{\alpha} \frac{1}{2} \sum_i \alpha_i x_i^T x_i - \sum_i \alpha_i y_i$$

$$\text{s.t. } \alpha_i y_i \geq 0 \quad \sum_i \alpha_i = 0$$

Is a quadratic programming problem
dual to the original one.

The solution $\hat{w} = \sum_i \alpha_i x_i$

\hat{w} determined by picking an x_i
s.t. $y_i(\hat{w}_0 + \hat{w}^T x_i) = 1$ $\Rightarrow \hat{w}_0 = y_i - \hat{w}^T x_i$

$\hat{y}(x) = \text{sgn}(\hat{w}_0 + \hat{w}^T x)$

Remember only certain α_i 's are non-zero. Turns out that when only a small fraction of training data becomes support vectors, the model has good generalization properties.

How do we generalize this linear maximum margin classifier to problems which require non linear decision surfaces.

The trick is to realize that we could replace $x_i^T x_j$ by a general positive definite kernel $K(x_i, x_j)$.

$$\min_{\alpha} \frac{1}{2} \sum_{i,j} \alpha_i K(x_i, x_j) \alpha_j - \sum_i \alpha_i y_i$$

$$\text{s.t. } \alpha_i \geq 0 \quad \sum_i \alpha_i = 0$$

$$\hat{y}(x) = \text{sgn}(\hat{w}_0 + \sum_i \alpha_i K(x_i, x))$$

With \hat{w}_0 determined from a support vector.

For example: one could use kernels like

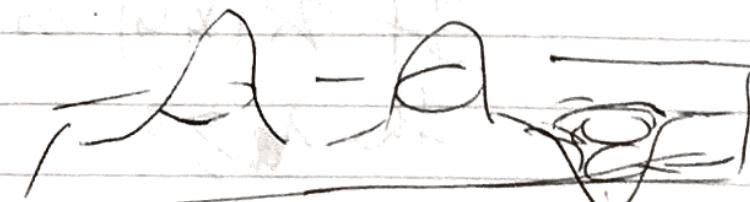
$$K(x, y) = (x \cdot y + 1)^P \rightarrow \text{Gaussian radial basis functions (RBF)}$$

$$K(x, y) = e^{-\|x - y\|^2 / 2\sigma^2}$$

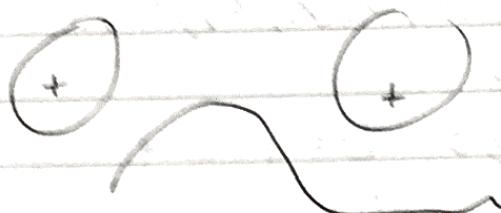
$$K(x, y) = \tanh(\pi x \cdot y - \delta)$$

Decision Surface RBF

$$\sum_i \alpha_i e^{-\|x_i - x\|^2 / 2\sigma^2}$$



Level sets of RBF kernel
could be as complex as possible
(+) (+) with enough data.



With a small
enough σ^2 , it
can approximate
any decision boundary.

Quite interestingly using $K(x, y)$
is equivalent to mapping x space
to an infinite dimensional feature
space by a nonlinear mapping
think of $K(x, y)$ as an infinite dimensional operator
 $\int K(x, y) \Phi(y) dy = \Phi(x)$

If K is symmetric we have infinitely
many eigenvalues $u_\alpha(x)$ with
 $\int K(x, y) u_\alpha(y) dy = \lambda_\alpha u_\alpha(x)$

and

$$K(x, y) = \sum_a \lambda_a u_a(x) u_a(y)$$

If $K(x, y)$ is positive definite,
 $\int K(x, y) g(x) g(y) dx dy > 0$
the $\lambda_a \geq 0$.

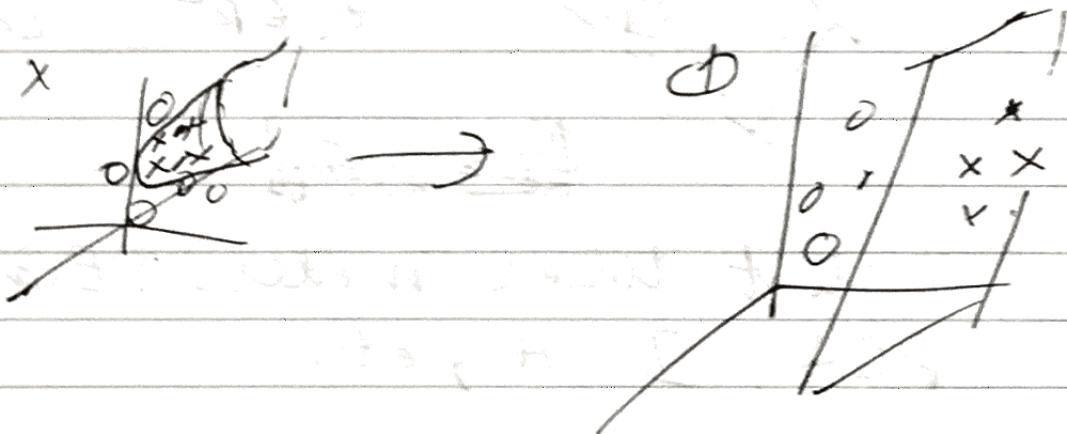
Define $\sqrt{x_i} \cdot u_i = \varphi_i$

$$K(x, y) = \sum_a \varphi_a(x) \varphi_a(y)$$

$$\text{So } K(x_i, x_j) = \sum_a \varphi_a(x_i) \varphi_a(x_j)$$

$$\text{This mean } K(x_i, x_j) = \Phi(x_i)^T \Phi(x_j)$$

$$\text{where } \Phi(x) = (\varphi_1(x), \varphi_2(x), \varphi_3(x), \dots)$$



Similar idea's could be used
for regression as well; $L_d(y, \hat{y}) \begin{cases} 0 & \text{if } y - \hat{y} \leq 0 \\ |y - \hat{y}| & \text{otherwise} \end{cases}$
 $y = w + w^T x, \min_{w, w^T} \sum_{i=1}^n L_d(y_i, \hat{y}_i)^2$

As an aside, exponential of a hinge loss
could be thought of as a mixture of many gaussians
with different scales, allowing a probability interpretation
of SVM's.

$$e^{-2a+} = \int_0^\infty e^{-\frac{1}{2} \frac{(a+\lambda)^2}{\lambda}} d\lambda$$

$$\text{Trick } \int e^{-\frac{1}{2} \left(\frac{a}{\lambda} + \lambda \right)^2} d\lambda = e^{-\frac{2a}{\lambda}} \int e^{-\frac{1}{2} \left(\frac{a}{\lambda} - \lambda \right)^2} d\lambda$$