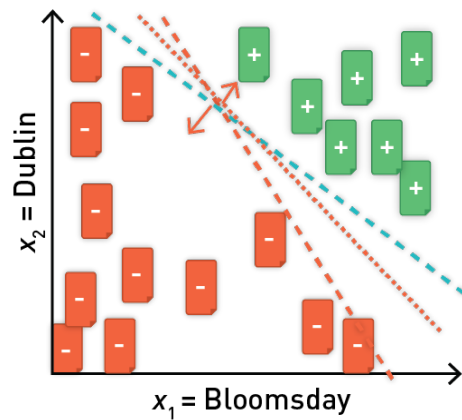


Linear Discriminant Model

Training Data

E.g.	x_1	y
1	3	-1
2		+1
...
L	0	-1

Many hyperplanes exist.



A linear discriminant function is linear in the components of X .

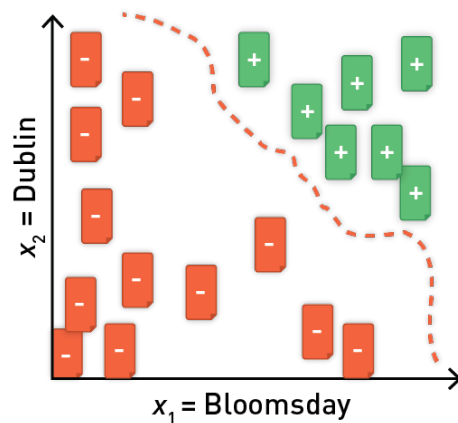
$$y = ax_1 + b = 0$$

$$f(X) = \begin{cases} -1 & \text{if } y < 0 \\ 0 & \text{if } y = 0 \\ 1 & \text{if } y > 0 \end{cases}$$

$$\text{Class}(X) = \text{sign}(\langle W, X \rangle + b)$$

Nonlinear Discriminant Model

Training Data			
E.g.	x_1	x_2	y
1	3	0	-1
2			+1
...
L	0	4	-1



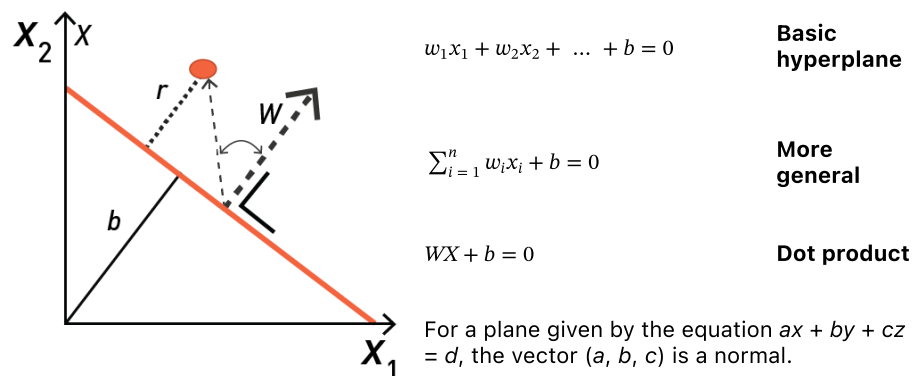
A nonlinear discriminant function is nonlinear in the components of X .

$$y = ax_1^2 + bx_2^3 + c = 0$$

$$f(X) = \begin{cases} -1 & \text{if } y < 0 \\ 0 & \text{if } y = 0 \\ 1 & \text{if } y > 0 \end{cases}$$

$$\text{Class}(X) = \text{sign}(\langle W, X \rangle + b)$$

Geometry: Linear Separators



Represent a hyperplane, H , in terms of vector W and scalar b .

- W determines the orientation of the hyperplane or discriminant plane.
- b denotes the offset (perpendicular distance) from plane to origin.

$$r = \frac{W^T X + b}{\sqrt{w_1^2 + \dots + w_n^2}}$$

Perpendicular distance from point X to a hyperplane

Empirical Risk Minimization

- This provides a criterion to decide on h .

$$\min_{h \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N \text{loss}(\mathbf{x}_i, \mathbf{y}_i; h)$$

- Background preferences over h can be included in regularized empirical risk minimization.

$$\min_{h \in \mathcal{H}} \frac{1}{N} \sum_{i=1}^N \text{loss}(\mathbf{x}_i, \mathbf{y}_i; h) + R(h)$$

Loss Terms: Zero-One and Hinge Loss

- Zero-one loss**, or L_{01} , counts the number of mistakes (gold standard).

$$L_{01}(m) = \begin{cases} 0 & \text{if } m \geq 0 \\ 1 & \text{if } m < 0 \end{cases}$$

- Hinge loss**, or L_{hinge} : Loss term for soft-margin SVM

$$\begin{aligned} J(w) &= \frac{1}{2} \|w\|^2 + \sum_i \max(0, 1 - y^i w^T x^i) \\ &= \frac{1}{2} \|w\|^2 + \sum_i \max(0, 1 - m_i(w)) \\ &= R_2(w) + \sum_i L_{\text{hinge}}(m_i) \end{aligned}$$

Loss Terms: Log Loss

- This is equivalent to the cross-entropy loss function used to train a logistic regression model.

$$J(w) = \lambda \|w\|^2 + \sum_i y^i \log g_w(x^{(i)}) + (1 - y^i)(\log 1 - g(x^{(i)})), y^i \in (0, 1)$$

- Simplify log conditional likelihood to get a more succinct loss component.

$$J(w) = \lambda \|w\|^2 + \sum_i \log 1 + e^{-y^{(i)} f_w(x^{(i)})}$$

- Where m is defined:

$$L(m) = \log 1 + e^{-m}$$

$$m^i = y^{(i)} f_w(x^{(i)})$$

$$y^{(i)} = \begin{cases} -1 & \text{if } y^{(i)} = 0 \\ 1 & \text{if } y^{(i)} = 1 \end{cases}$$

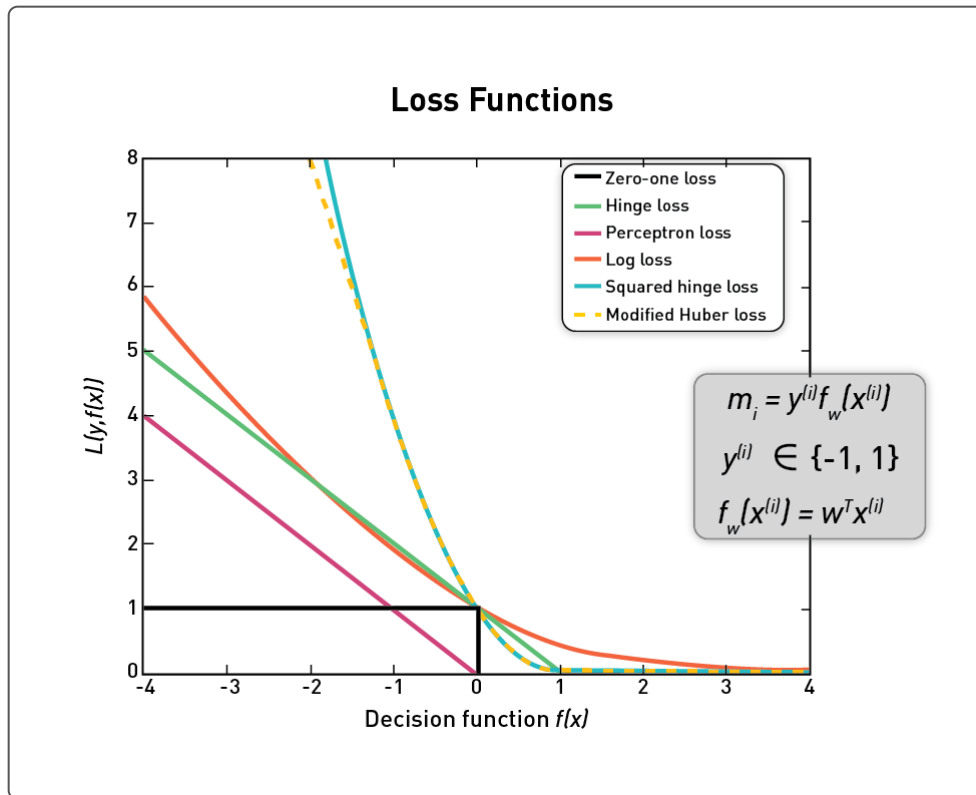
Other Loss Terms

- **Squared-error:** For linear regression

$$L_2(m) = (f_w(x) - y)^2 = (m - 1)^2$$

- **Boosting:** A greedy optimization of the exponential loss term

$$J(w) = \lambda R(w) + \sum_i \exp(-y^{(i)} f_w(x^{(i)}))$$
$$L_{exp}(m_i) = \exp(-m_i(w))$$



Loss Functions: A Unifying View

Loss function consists of:

- **Loss term** $Lm_i w$, expressed in terms of the margin of each training example
- **Regularization term** $\lambda R w$, expressed as a function of the model complexity

$$J(w) = \sum_i \underbrace{L(m_i(w))}_{\text{Loss term}} + \underbrace{\lambda R(w)}_{\text{Regularization term}}$$

$$m_i = y^{(i)} f_w(x^{(i)})$$

$$y^{(i)} \in \{-1, 1\}$$

$$f_w(x^{(i)}) = w^T x^{(i)}$$

Machine Learning Objectives

Objective function:

$$J(w) = \sum_i \underbrace{L(m_i(w))}_{\text{Loss term}} + \underbrace{\lambda R(w)}_{\text{Regularization term}}$$

Almost all machine learning objectives are optimized using this update:

$$w \leftarrow w - \alpha \cdot \sum_{i=1}^n g(J(w)x_i, y_i)$$

Update weight vector with gradient of the objective function.

- w is a vector of dimension d .
- We're trying to find the best w via optimization.

Distributed Gradient Descent

- Ordinary least squares (OLS)
- Logistic regression
- Bayesian logistic regression
- Perceptron
- SVMs and their many learning algorithms

Distributed Gradient Descent: Linear Regression Example

Master-Slave Process

- Initialize model parameters; assume a weight vector $W = (0, 0, \dots, 0)$; Gradient $= (0, 0, \dots, 0)$
- While not converged
 - MASTER: Broadcast model (i.e, weight vector) to the worker nodes
 - MASTER launches MapReduce jobs
 - Mappers (*many* mappers) to compute partial gradients over the respective training data subsets (chunks)

Init $g = (0, 0, \dots, 0)$.

For each training example:

 Compute partial gradient for each training example.

 Combine in memory: $g = \sum_i gW; X_i, Y_i$.

Finally yield the partial gradient g .
 - Reducer (single reducer)

Initialize full gradient: $G = (0, 0, \dots, 0)$.

For each partial gradient g :

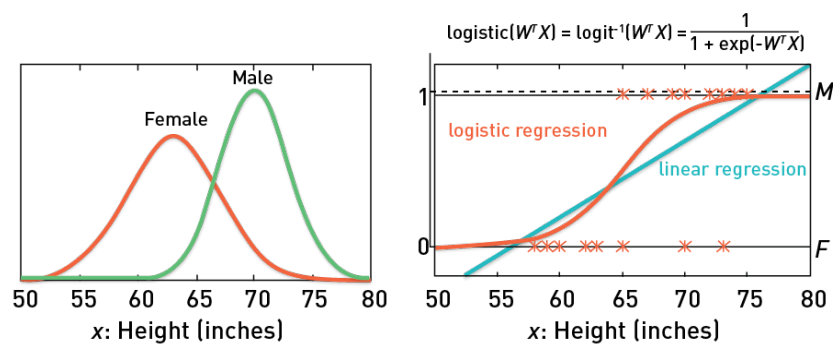
 Aggregate partial gradients: $g = \sum_m g_m$.

Yield full gradient G .
- MASTER $W = W + \alpha G$ // update the weight vector
- End-While

Linear regression
Goal: Learn a weight vector W .
Gradient is defined here:
 $gW; X_i, Y_i = y^i - W_i X^i$

Logistic Regression in One Slide

Example: Predict the gender ($y = M/F$) of a person given their height ($x = \text{a number}$).



$$py = Mx = \beta_0 + \beta_x \text{ (linear regression)}$$

$$\ln \frac{py = Mx}{py = Fx} = \beta_0 + \beta_x \text{ (logistic regression)}$$

Assumptions

- $y \in \{-1, +1\}$, x is feature vector, and p is defined as:

$$p = \Pr(y = +1 \mid x)$$

- Linearly related to features after logit transformation:

$$\log \frac{p}{1-p} = w^T \cdot x - b$$

- Then probability is a logistic function of $\mathbf{w} \cdot \mathbf{x}$:

$$p = \frac{1}{1 + \exp -w^T \cdot x + b}$$

- Where:
 - w is the coefficient vector.
 - b is the intercept.

Maximum Likelihood Expectation (MLE)

Maximize the log likelihood:

$$I(W) = \ln \prod_i P_i$$

$$= \ln \prod_i \left(\frac{1}{1 + \exp(-\mathbf{w}^T \cdot \mathbf{x}_i + b)} \right)^{\frac{1+y_i}{2}} \left(1 - \frac{1}{1 + \exp(-\mathbf{w}^T \cdot \mathbf{x}_i + b)} \right)^{\frac{1-y_i}{2}}$$

Maximizing this log likelihood is equal to minimizing the following:

$$I(W) = \sum_i \log(1 + \exp(-y\mathbf{w}^T \mathbf{x}_i))$$

$$L(m) = \log 1 + e^{-m}$$

$$m^i = y^{-(i)} f_w(\mathbf{x}^{(i)})$$

- Minimize the **NEG** log joint conditional likelihood.
- The conditional likelihood θ given data x and y is $L(\theta; y|x) = p(y|x; \theta)$.

Estimating Parameters Using Gradient Descent

- No closed-form solution to maximizing $l(W)$ with respect to W
- One common approach: Use gradient descent, working with the gradient, which is the vector of partial derivatives

Estimating Parameters Using Gradient Descent (cont.)

The i th component of the vector gradient has the form:

$$I(W) = \sum_i Y^i \left(w_0 + \sum_i^n w_i X_i^i \right) - \ln \left(1 + \exp \left(w_0 + \sum_i^n w_i X_i^i \right) \right)$$

$$\frac{\partial I(W)}{\partial w_i} = \sum_l X_i^l \left(Y^l - \hat{P}(Y^l = 1 | X^l, W) \right)$$

$$w_i \leftarrow w_i + \eta \sum_l X_i^l \left(Y^l - \hat{P}(Y^l = 1 | X^l, W) \right)$$

- Beginning with initial weights of zero, we repeatedly update the weights in the direction of the gradient, changing the i th weight according to this formula, where η is a small constant (e.g., 0.01), which determines the step size.
- Effectively, we are pulling the weight vector closer to the examples where we make mistakes.

Gradient Descent

- Objective function (logloss):

$$\sum_i \log (1 + \exp (-y (\mathbf{w}^T \mathbf{x}_i + b)))$$

- Gradient:

$$\nabla \mathbf{w} = \sum_i -y \left(1 - \frac{1}{1 + \exp(-y(\mathbf{w}^T \mathbf{x}_i + b))} \right) \cdot \mathbf{x}_i$$

- Update \mathbf{w} until it converges:

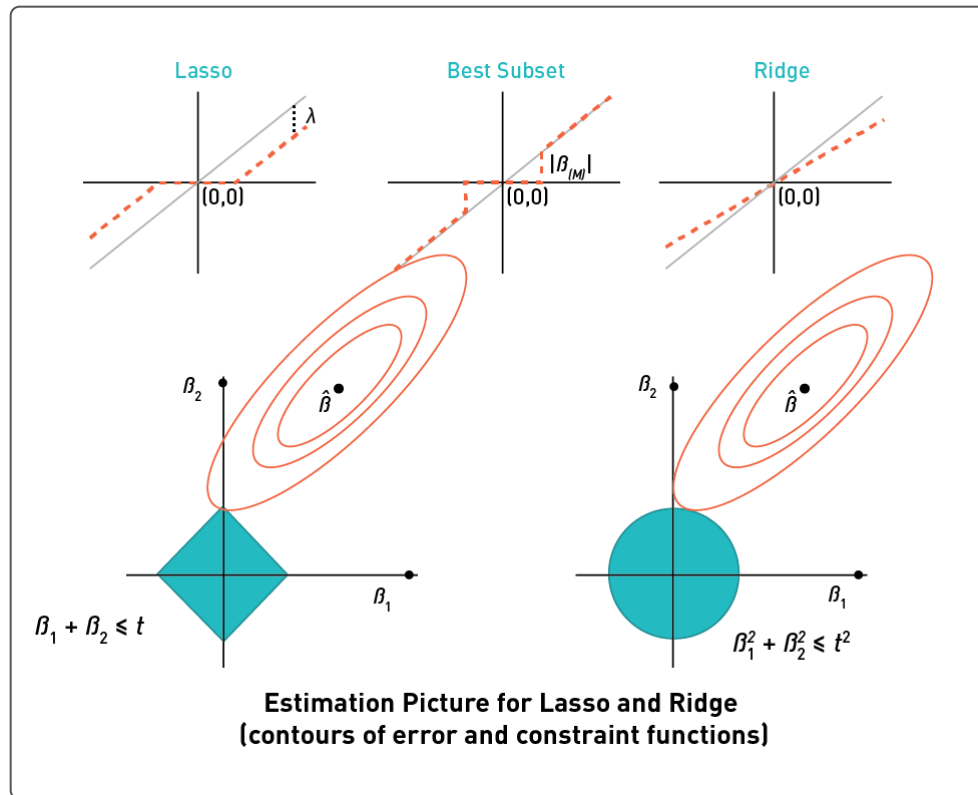
$$\mathbf{w} = \mathbf{w} - \eta \cdot \nabla \mathbf{w}$$

Regularization

Shrinkage of \mathbf{w}

- **Lasso:** L1 norm regularization, $|\mathbf{w}|$
- **Ridge:** L2 norm regularization, \mathbf{w}^2

Add regularization to prevent overfitting.



Gradient Descent: Lasso

- Objective function (logloss and lasso):

$$\sum_i \log(1 + \exp(-y(\mathbf{w}^T \mathbf{x}_i + b))) + \lambda |\mathbf{w}|$$

- Gradient:

$$\nabla \mathbf{w} = \sum_i -y \left(1 - \frac{1}{1 + \exp(-y(\mathbf{w}^T \mathbf{x}_i + b))} \right) \cdot \mathbf{x}_i + \lambda (\mathbf{1}_{>0}(\mathbf{w}) \cdot 2 - 1)$$

- Update w until it converges:

$$\mathbf{w} = \mathbf{w} - \eta \cdot \nabla \mathbf{w}$$

Gradient Descent: Ridge

- Objective function (logloss and ridge):

$$\sum_i \log(1 + \exp(-y(\mathbf{w}^T \mathbf{x}_i + b))) + \lambda \mathbf{w}^2$$

- Gradient:

$$\nabla \mathbf{w} = \sum_i -y \left(1 - \frac{1}{1 + \exp(-y(\mathbf{w}^T \mathbf{x}_i + b))} \right) \cdot \mathbf{x}_i + \lambda \mathbf{w}$$

- Update w until it converges:

$$\mathbf{w} = \mathbf{w} - \eta \cdot \nabla \mathbf{w}$$

Python Notebooks for Logistic Regression

- Notebook for logistic regression on a single node:
<http://nbviewer.ipython.org/urls/dl.dropbox.com/s/kurbw695jdvxib0/LogisticRegression-Single-Core-Notebook-.ipynb>
- Notebook for distributed logistic regression on Spark:
<http://nbviewer.ipython.org/urls/dl.dropbox.com/s/r20ff7q0yni5kiu/LogisticRegression-Spark-Notebook.ipynb>

Logistic Regression: Single-Node Code

Python—Gradient Descent

Input:

- data: feature information
- y: label information
- eta: learning rate
- iter_num: maximum iteration number
- regPara: regularization parameter
- stopCriteria: stop criteria

Output

- w: coefficients of boundary

[LogisticRegression-Single-Core-Notebook-.ipynb](#)

Note: This is the single-node implementation of logistic regression. The distributed version will follow momentarily.

Please click on the link and take some time out to review this notebook.

```
# gradient descent (and with NO stochasticity!)
# Objective Function
# minw λ/2 w'w + 1/m ∑ log(1+exp(yi(w'xi - b)))
# gradient
# -y*(1-1/(1+exp(yi(w'xi - b))))*x

def
logisticReg_GD(data,y,w=None,eta=0.05,iter_num=500,regPara=0.01,stopCriteria=
0.0001,reg="Ridge"):
    data = np.append(data,np.ones((data.shape[0],1)),axis=1)
    if w is None:
        w = np.random.normal(size=data.shape[1])
    for i in range(iter_num):
        wxy = np.dot(data,w)*y
        g = np.dot(data.T,-y*(1-1/(1+np.exp(-wxy))))/data.shape[0]
        #Gradient of log loss
        wreg = w*1
        wreg[-1] = 0 #last value of weight vector is bias term;
        ignore in regularization
        if reg == "Ridge":
            wreg = w*1
            wreg[-1] = 0 #last value of weight vector is bias term;
            ignore in regularization
        elif reg == "Lasso":
            wreg = w*1
            wreg[-1] = 0 #last value of weight vector is bias term;
            ignore in regularization
            wreg = (wreg>0).astype(int)*2-1
        else:
            wreg = np.zeros(w.shape[0])
        wdelta = eta*(g+regPara*wreg) #gradient: log loss + regularized term
        if sum(abs(wdelta))<=stopCriteria*sum(abs(w)): #Convergence condition
            break
        w = w - wdelta
    return w
```

Distributed Logistic Regression: PySpark Code

[LogisticRegression-Spark-Notebook.ipynb](#)

```
# gradient descent (and with NO stochasticity!)
# Objective Function
# minw λ/2 w'w + 1/m ∑i log(1+exp(yi(w'xi - b)))
# gradient
# -y*(1-1/(1+exp(yi(w'xi - b))))*x

def logisticReg_GD_Spark(data,y,w=None,eta=0.05,iter_num=500,regPara=0.01,
stopCriteria=0.0001,reg="Ridge"):
    #eta learning rate
    #regPara
    dataRDD = sc.parallelize(np.append(y[:,None],data,axis=1)).cache()
    if w is None:
        w = np.random.normal(size=data.shape[1]+1)
    for i in range(iter_num):
        w_broadcast = sc.broadcast(w)
        g = dataRDD.map(lambda x: -x[0]*(1-1/(1+np.exp(-x[0]
        *np.dot(w_broadcast.value,np.append(x[1:],1)))))) \
        *np.append(x[1:],1)).reduce(lambda x,y:x+y)/data.shape[0]
        # Gradient of logloss
        if reg == "Ridge":
            wreg = w*1
            wreg[-1] = 0 #last value of weight vector is bias term;
            ignore in regularization
        elif reg == "Lasso":
            wreg = w*1
            wreg[-1] = 0 #last value of weight vector is bias term;
            ignore in regularization
            wreg = (wreg>0).astype(int)*2-1
        else:
            wreg = np.zeros(w.shape[0])
        wdelta = eta*(g+regPara*wreg) #gradient: hinge loss + regularized term
        if sum(abs(wdelta))<=stopCriteria*sum(abs(w)): # converged as updates
            to weight vector are small
            break
        w = w - wdelta
    return w
```

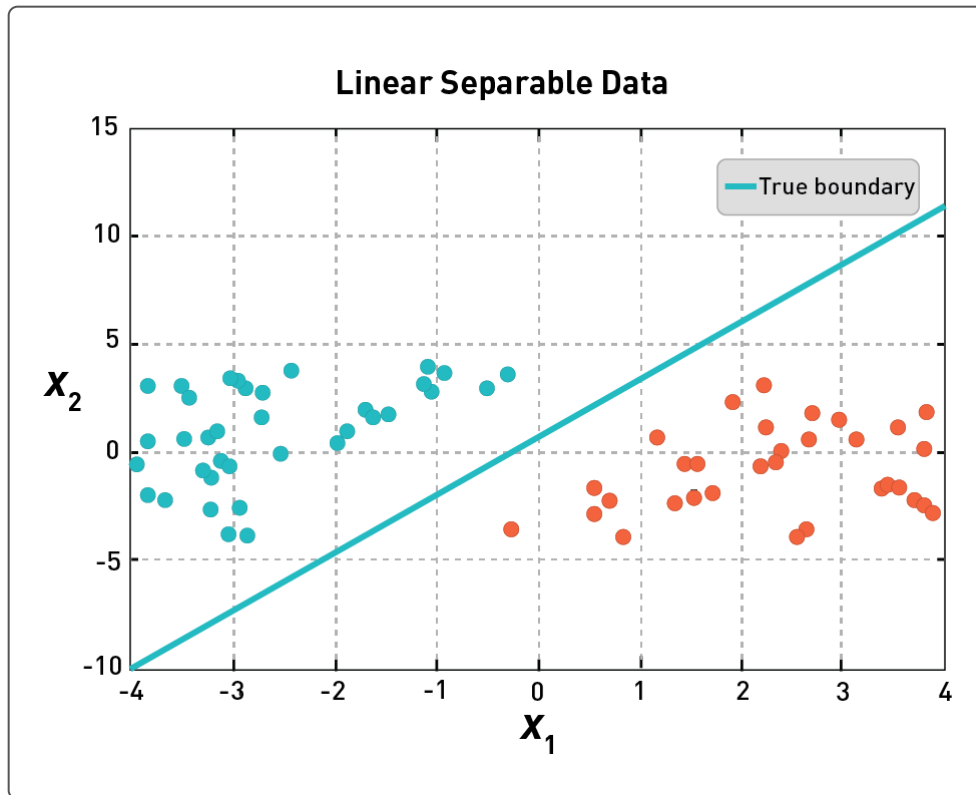
How to Run the Code

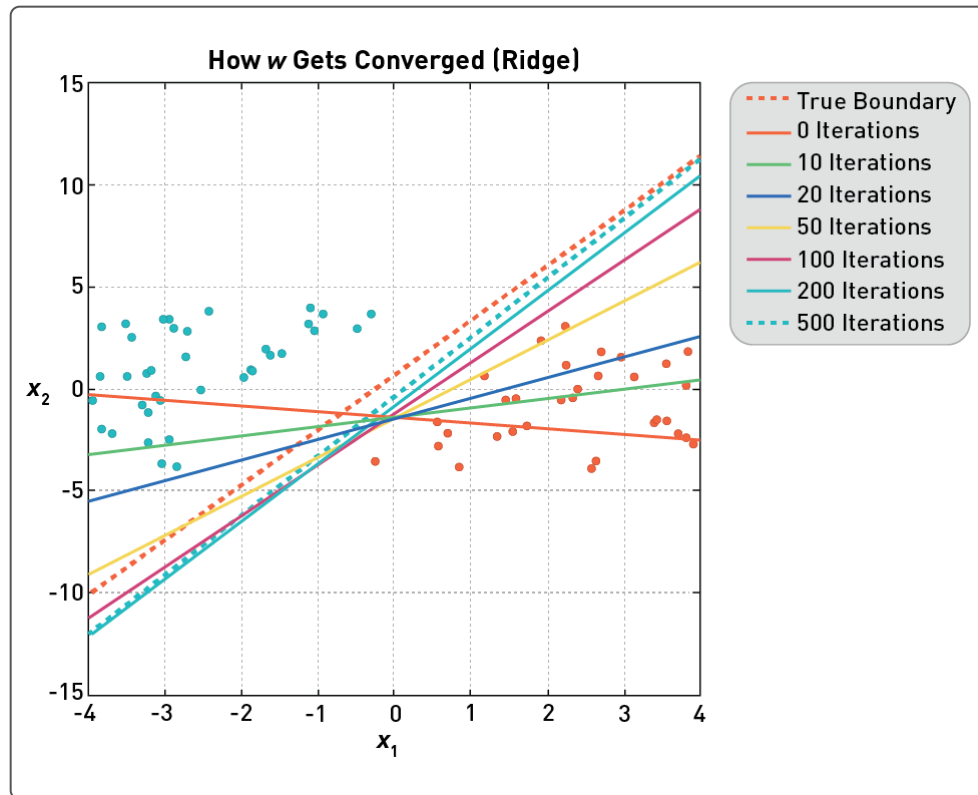
- **Ridge logistic regression**

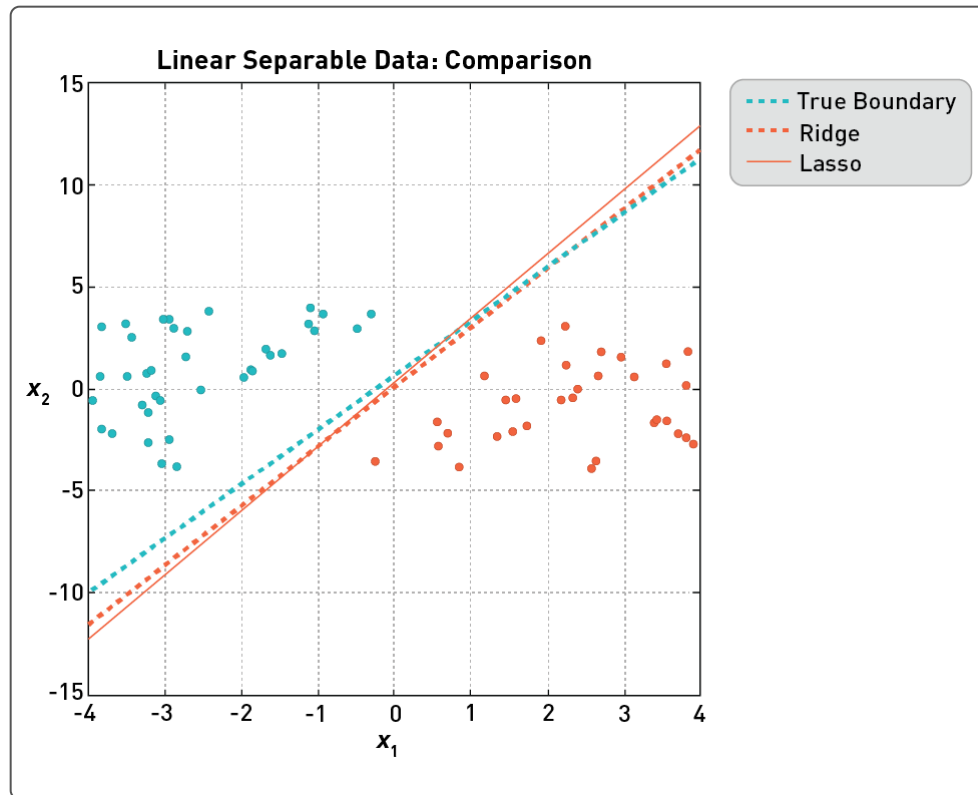
- `np.random.seed(400)`
`w_gd_spark_sep_ridge =`
`logisticReg_GD_Spark(data_sep, y_sep, reg="Ridge")`
`print w_gd_spark_sep_ridge`
 - `[1.84106359 -0.63166074 0.00472968]`

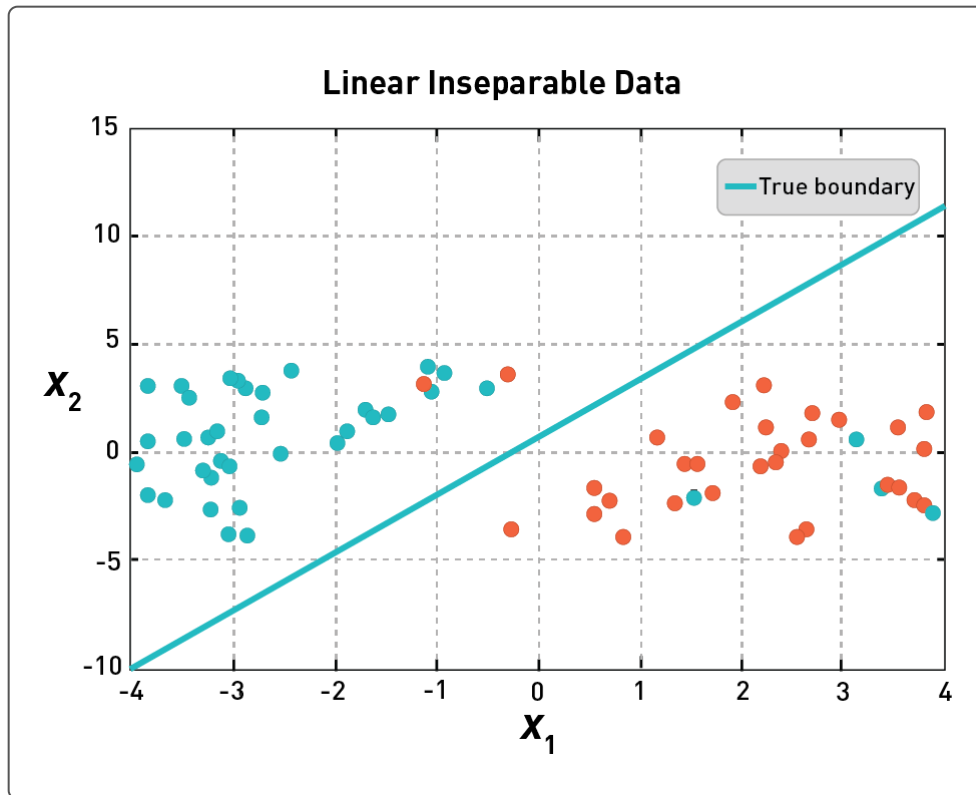
- **Lasso logistic regression**

- `np.random.seed(400)`
`w_gd_spark_sep_lasso =`
`logisticReg_GD_Spark(data_sep, y_sep, reg="Lasso")`
`print w_gd_spark_sep_lasso`
 - `[1.90212754 -0.60763246 0.18656826]`









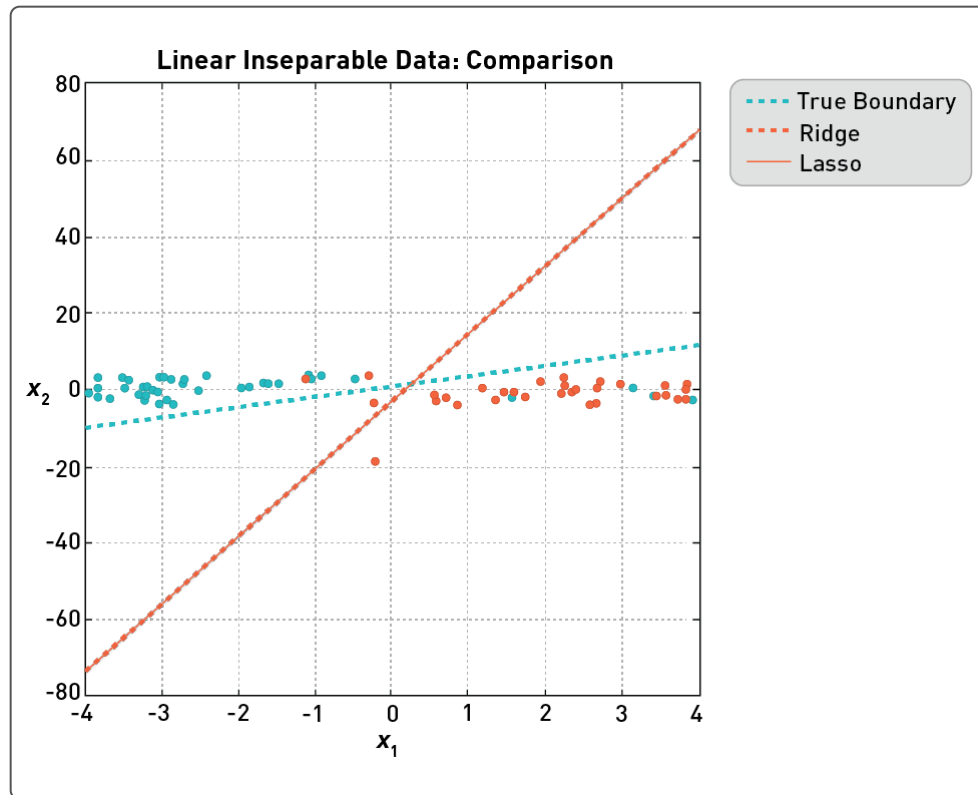
Linear Inseparable Data: Code

- **Ridge logistic regression**

- `np.random.seed(400)`
 - `w_gd_insep_ridge = logisticReg_GD(data_insep, y_insep, reg="Ridge")`
 - `print w_gd_insep_ridge`
 - `[0.88960208 -0.06413379 -0.24526051]`

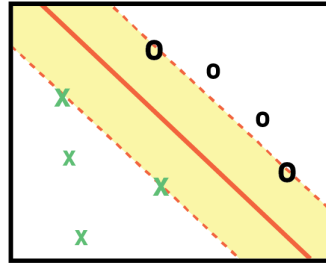
- **Lasso logistic regression**

- `np.random.seed(400)`
 - `w_gd_insep_lasso = logisticReg_GD(data_insep, y_insep, reg="Lasso")`
 - `print w_gd_insep_lasso`
 - `[0.88745901 -0.0500822 -0.14752302]`



Find Shortest Weight Vector

- Finding a separating hyperplane with the largest margin is equivalent to finding such a hyperplane with minimal W .
- Solve using optimization approaches.



$$\text{Max } \frac{1}{\|w\|^2}$$

Subject to constraints

$$y_i w \cdot x_i + b - 1 \geq 0 \quad \forall i = 1, \dots, L$$

or

$$\text{Min } \frac{\|w\|^2}{2}$$

Subject to constraints

$$y_i w \cdot x_i + b - 1 \geq 0 \quad \forall i = 1, \dots, L$$

Hard SVMs vs. Soft SVMs

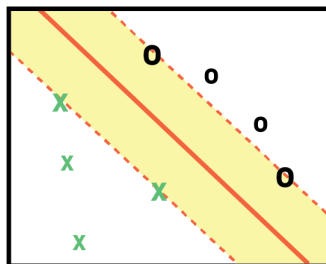
Hard SVM learning algorithm

Must classify each example correctly

$$\min Wb = 0.5 \times \|w\|^2$$

Subject to:

$$y_i W X_i + b \geq 1 \quad \forall i = 1, \dots, L$$



Soft SVM learning algorithm

Has trade-off between margin and error

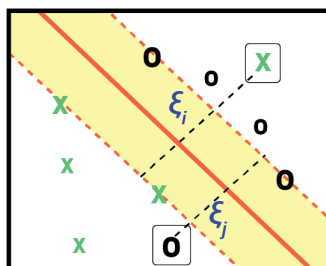
Associates a slack variable ξ_i with each example

$$\min Wb = 0.5 \times \|w\|^2 + C \sum_{i=1}^L \xi_i$$

Subject to:

$$y_i W X_i + b + \xi_i \geq 1 \quad \forall i = 1, \dots, L$$

and $\xi_i \geq 0 \quad \forall i = 1, \dots, L$



Kernels

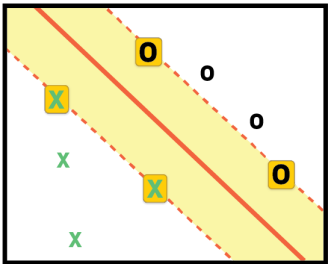
- Linear
- Polynomial
- Radial basis function

Primal and Dual Learning

- **Primal learning** (e.g., perceptron) involves learning weight values associated with each term or feature.
- **Dual learning** involves learning weight values associated with each example so that the margin, λ , is maximum.

Dual Learning

Learning weight values associated with each example so that the margin, λ , is maximum



α	Wgt Vector		w_0	w_1	w_n	y
	Instance/Attr		x_0	x_1	x_2	...	x_n	
1.2	1		1	3	0	...	7	-1
0	2		1			...		+1
\vdots	\vdots		\vdots	\vdots	\vdots	...	\vdots	\vdots
4.2	L		1	0	4	...	8	-1

Dual Learning (cont.)

The dual problem is solved for multipliers, not for w and b .

- W is then obtained as a linear combination of data vectors, by the constraint:
 - $w - \sum_i \alpha_i y_i x_i = 0$
- b can be found by complementary Karush-Kuhn-Tucker (KKT) condition:
 - $\alpha_i = 0 \Rightarrow y_i f(\vec{x}_i) > 1$
 - $\alpha_i > 0 \Rightarrow y_i f(\vec{x}_i) = 1$ **SV (active constraint)**

$$W = \sum_{i=1}^L \alpha_i y_i X_i$$

and $b = y_{SV} - \langle W, X_{SV} \rangle$

$$\langle W, X \rangle = \left\langle \left(\sum_{i=1}^L \alpha_i y_i X_i \right), X \right\rangle = \sum_{i=1}^L \alpha_i y_i \langle X_i, X \rangle$$

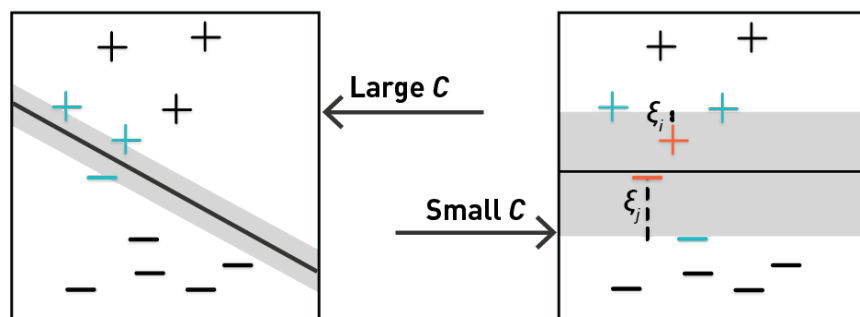
Soft Margin SVM (Primal)

$$\min Wb = 0.5 \times \|w\|^2 + C \sum_{i=1}^L \xi_i$$

Subject to:

$$y_i W X_i + b + \xi_i \geq 1 \quad \forall i = 1, \dots, L$$

$$\text{and } \xi_i \geq 0 \quad \forall i = 1, \dots, L$$

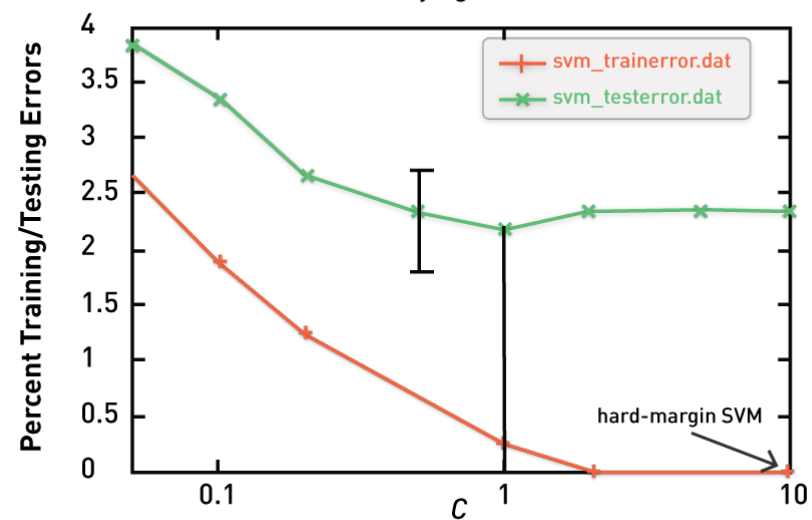


$C > 0$ is a *regularization parameter*. It takes the form of a tuning constant that controls the size of the slack variables and balances the two terms in the minimizing function.

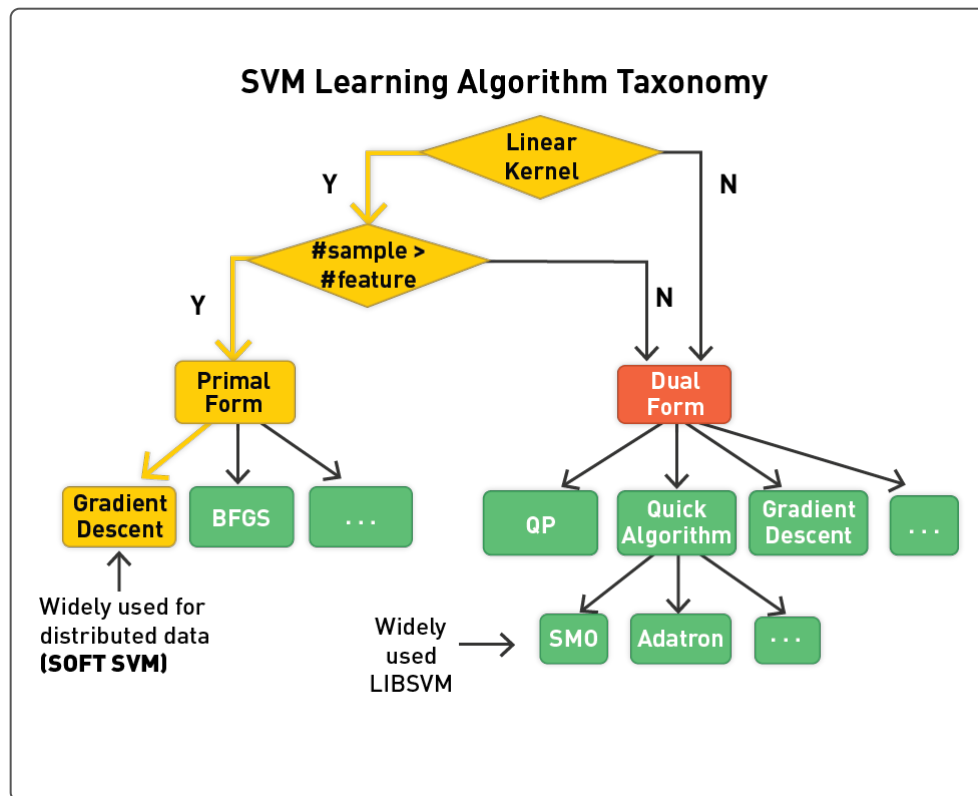
- Large $C \rightarrow$ hard margin (allows very few errors)
- Small $C \rightarrow$ allows a lot of slack and therefore a large margin

Controlling Soft-Margin Separation

Example: Reuters "acq"
Varying C



Observation: Typically no local optima, but not necessarily



Soft SVM via Unconstrained Optimization

Soft margin SVM (primal)

$$\min Wb = 0.5 \times \|w\|^2 + C \sum_{i=1}^L \xi_i$$

- Subject to:

$$y_i W X_i + b + \xi_i \geq 1 \quad \forall i = 1, \dots, L$$

$$\xi_i \geq 0 \quad \forall i = 1, \dots, L$$

Objective function

- Min $\frac{\|w\|^2}{2}$
- Subject to constraints: $y_i w \cdot x_i + b - 1 \geq 0 \quad \forall i = 1, \dots, L$

Soft SVM via Unconstrained Optimization (cont.)

Lagrangian function

$$L_p(\mathbf{w}) = \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle - \sum \alpha_i [y_i (\langle \mathbf{w}, x_i \rangle + b) - 1]$$

$$\alpha \geq 0$$

SVM Learning as Optimization Problem

Soft margin SVM (primal)

$$\min_w \quad \frac{1}{2} \|w\|^2 + C \sum_{i=1}^L \xi_i$$

- Subject to:

$$y_i w'x_i + b + \xi_i \geq 1 \quad \forall i = 1, \dots, L$$

$$\xi_i \geq 0 \quad \forall i = 1, \dots, L$$

Regularized hinge loss

$$\min_w \quad \frac{\lambda}{2} w'w + \frac{1}{m} \sum_{i=1}^m \max(0, 1 - y_i w'x_i)$$

Trade-off between margin and loss Size of the margin Expected hinge loss on the training set Positive for correctly classified examples, otherwise negative

SVM Learning as Optimization Problem (cont.)

Regularized hinge loss

$$\min_w \frac{\lambda}{2} w'w + \frac{1}{m} \sum_i (1 - y_i (w'x_i - b))_+$$

- $1 - z_+ = \max(0, 1 - z)$ (hinge loss)
- First summand is a quadratic function; the sum is a piecewise linear function
- The whole objective: Piecewise quadratic

Gradient of Regularized Hinge Loss

Regularized hinge loss

$$\min_w \frac{\lambda}{2} w'w + \frac{1}{m} \sum_i (1 - y_i (w'x_i - b))_+$$

Gradient of regularized hinge loss

- Hinge loss and regularization loss:

$$\begin{array}{ll} \lambda w & \text{if } y_i w'x_i - b > 1 \\ \lambda w + y_i x_i & \text{Otherwise} \end{array}$$

- $W_{t+1} = w_t + \text{average}(\text{gradient})$
- $W^{t+1} = w_t + \text{average}(\text{regularization} + \text{hinge loss})$

Distributed Gradient Descent: Linear Soft SVMs

Master-Slave Process

- Initialize model parameters; assume a weight vector $W = (0, 0, \dots, 0)$; Gradient = $(0, 0, \dots, 0)$
- While not converged
 - MASTER: Broadcast model (i.e, weight vector) to the worker nodes
 - MASTER launches MapReduce jobs
 - Mappers (*many* mappers) to compute partial gradients over the respective training data subsets (chunks)

Init $g = (0, 0, \dots, 0)$.
Combine in memory:
 $g = \sum_i gW; X_i, Y_i$.

Gradient of regularized hinge loss: λw if $y_i w' x_i - b > 1$
 $\lambda w + y_i x_i$ Otherwise

Finally yield the partial gradient g .
 - Reducer (single reducer)

Initialize full gradient: $G = (0, 0, \dots, 0)$.

For each partial gradient g

Aggregate partial gradients: $g = \sum_m g_m$.

Yield full gradient G .
 - MASTER $W = W + \alpha G$ // update the weight vector
- End-While

PySpark: SVM via Gradient Descent (Single Node)

Link to notebook: [SVM-Notebook-Linear-Kernel-2015-06-19.ipynb](#)

```

1 #gradient descent (and with no stochasticity!)
2 #Objective Function
3 #minw  λ/2  w'w  +  1/m ∑i(1 - yi(w'xi - b))+
4 #gradient
5 #  λw          if yi(w'xi - b) > 1  #correctly classified
6 #  λw + yi xi  Otherwise            #incorrectly classified
7 def SVM_GD(data,y,w=None,eta=0.01,iter_num=1000,regPara=0.01,
  stopCriteria=0.0001):
8     data = np.append(data,np.ones((data.shape[0],1)),axis=1)
9     if w==None:
10        w = np.random.normal(size=data.shape[1])
11    for i in range(iter_num):
12        wxy = np.dot(data,w)*y #labeled margin
13        xy = -data*y[:,None]
14        zipv = zip(xy,wxy)
15        #Gradient of hinge loss: if wxy<0 then hinge loss is xy
16        g = sum((u for u,v in zipv if, v <1))/data.shape[0]
17        wreg = w*1 #weight vector
18        wreg[-1] = 0
19        #wreg = np.array([w[0],w[1],0])
20        wdelta = eta*(g+regPara*wreg)
21        if sum(abs(wdelta)) <= stopCriteria*sum(abs(w)):
22            break
23        w = w - wdelta
24    return w

```

PySpark: Distributed SVM

Link to notebook: [SVM-Notebook-Linear-Kernel-2015-06-19.ipynb](#)

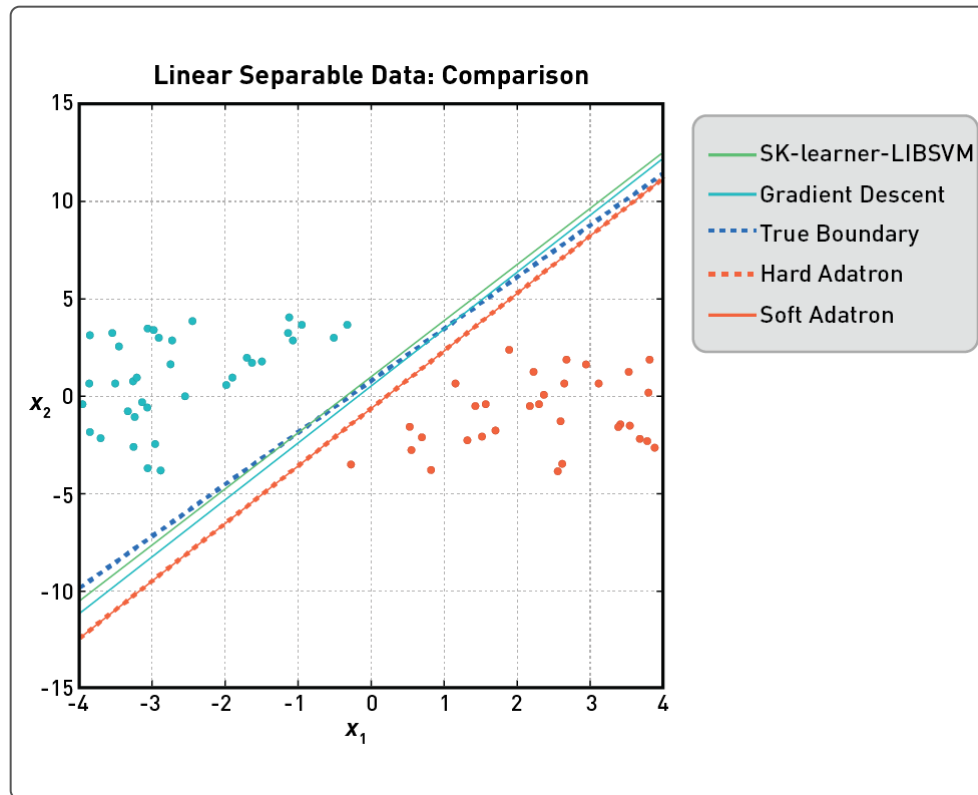
```

1 #gradient descent (and with no stochasticity!)
2 #Objective Function
3 #minw  $\lambda/2 \quad w'w + 1/m \sum_i (1 - y_i(w'x_i - b))^+$ 
4 #gradient
5 #  $\lambda w$  if  $y_i(w'x_i - b) > 1$  #correctly classified
6 #  $\lambda w + y_i x_i$  Otherwise #incorrectly classified
7
8 def SVM_GD_SPARK(data,y,w=None,eta=0.01,iter_num=1000,regPara=0.01,
  stopCriteria=0.0001):
9     #eta learning rate
10    #regPara
11    dataRDD = sc.parallelize(np.append(y[:,None],data,axis=1)).cache()
12    #prepend y to x
13    if w==None:
14        w = np.random.normal(size=data.shape[1]+1)
15    for i in range(iter_num): #label*margin
16        sv = dataRDD.filter(lambda x:x[0]*np.dot(w,np.append(x[1:],1))<1)
17        #Support vector? with label*margin<1
18        if sv.isEmpty(): #converged as no more updates possible
19            break #hinge loss component of gradient y*x and sum up
20        g = -sv.map(lambda x:x[0]*np.append(x[1:],1)).reduce(lambda
21        x,y:x+y)/data.shape[0] #gradient: total hinge l
22        wreg = w*1 #temp copy of weight vector
23        wreg[-1] = 0 #last value of weight vector is bias term;
24        ignore in regularization
25        wdelta = eta*(g+regPara*wreg) #gradient: hinge loss + regularized term
26        if sum(abs(wdelta))<=stopCriteria*sum(abs(w)):
27            #converged as updates to weight vector are small
28            break
29        w = w - wdelta
30    return w

```

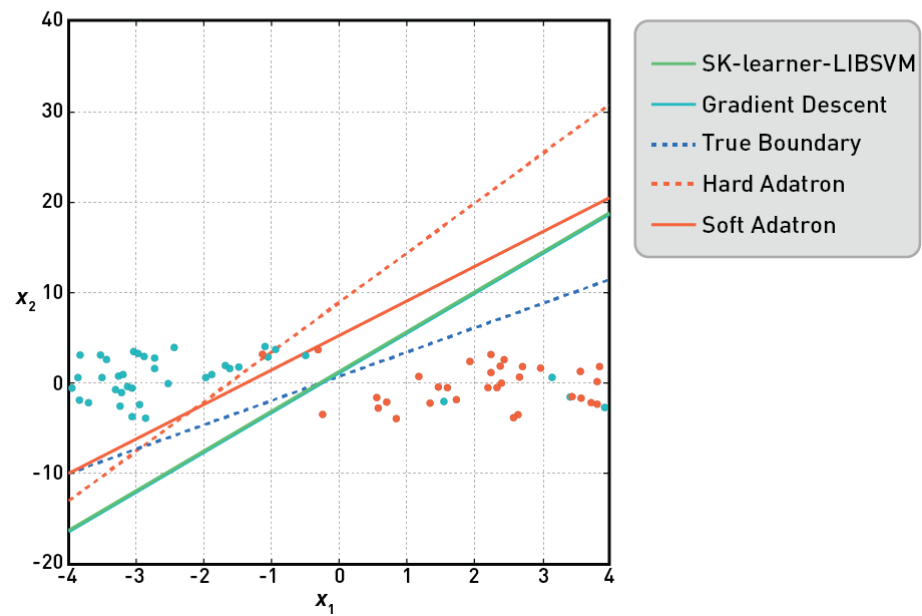
Linearly Separable Data: Code

```
• np.randomseed(400)
  w_gd_spark_sep = SVM_GD_SPARK(data_sep, y_sep)
  print w_gd_spark_sep
• array([ 0.93510664, -0.31799796,  0.10234458])
```



Linear Inseparable Data: Comparison

```
: comparison_plot(data_insep,y_insep,w_insep,w_sklearn_libsvm_insep,w_gd_insep,w_hard_adatron_insep,  
w_soft_adatron_insep,sep=False)
```



Supervised Machine Learning

- Loss functions
- General framework for gradient descent notebook
- Logistic regression at scale
- Perceptron
- SVMs