

INTRODUCTION TO MACHINE LEARNING COMPSCI 4ML3

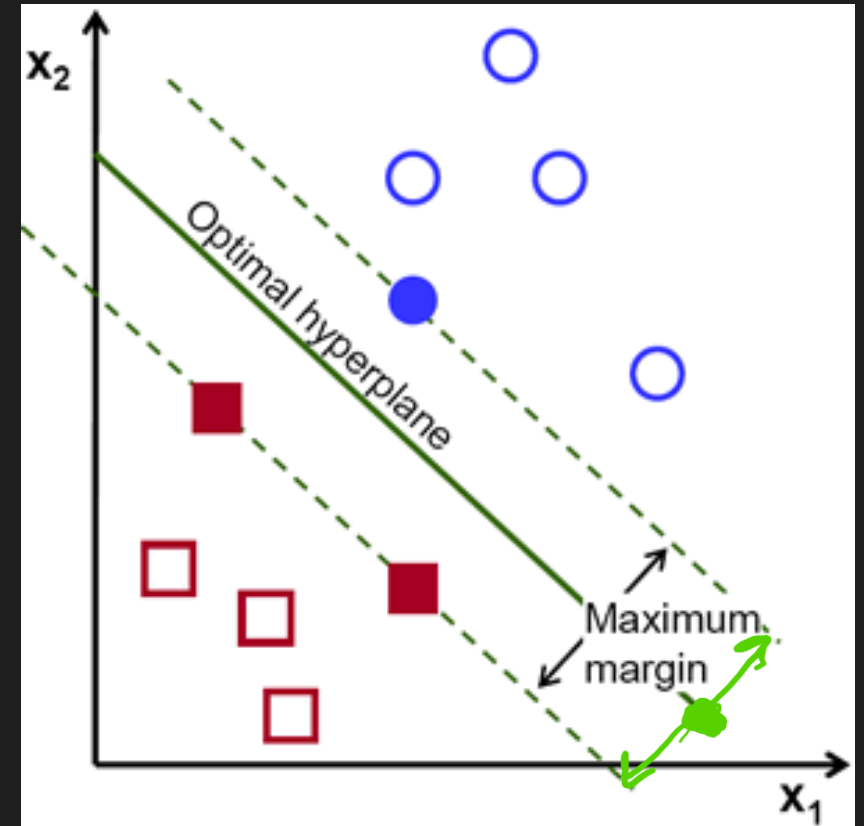
LECTURE 16

HASSAN ASHTIANI



SUPPORT VECTOR MACHINES

- PICK THE LINEAR SEPARATOR THAT MAXIMIZES THE “MARGIN”
- MORE ROBUST TO “PERTURBATION”
- LESS OVERFITTING PROBLEM!
 - WORKS WELL FOR HIGH-DIMENSIONAL DATA (?)
 - MORE ON THAT LATER!



DISTANCE OF A POINT TO A HYPERPLANE

THE EUCLIDEAN DISTANCE BETWEEN A POINT x AND THE HYPERPLANE PARAMETRIZED BY W IS (WHY?)

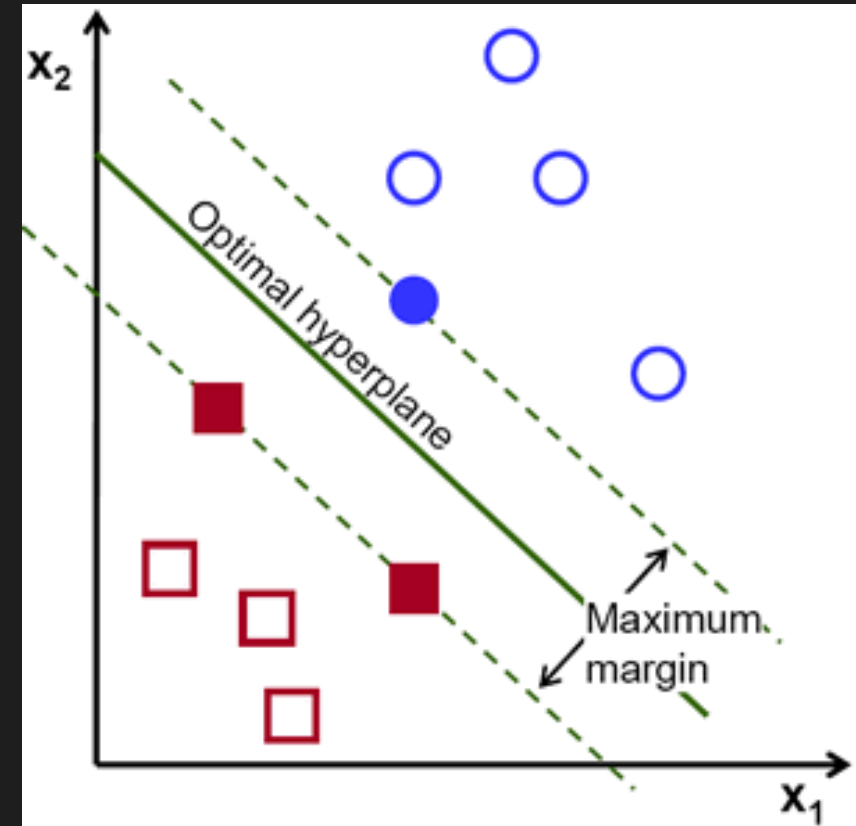
$$\frac{|W^T x + b|}{\|W\|_2}$$

- ONLY THE DIRECTION OF W MATTERS (NOT $\|W\|_2$)
- ASSUME $\|W\|_2=1$, THEN THE DISTANCE IS

$$|W^T x + b|$$

MAXIMUM MARGIN HYPERPLANE

- LET THE HYPERPLANE BE PARAMETRIZED BY W AND b
- ASSUME $\|W\|_2 = 1$
- (W, b) HAS A γ MARGIN IF
 - $W^T x + b > \gamma$ FOR EVERY BLUE x , AND
 - $W^T x + b < -\gamma$ FOR EVERY RED x





THE VERSION WITH “BIAS”

Hard-SVM

$$w' \approx \frac{w}{\|w\|} \quad \text{such that } |w' \cdot x_i| \geq 1$$

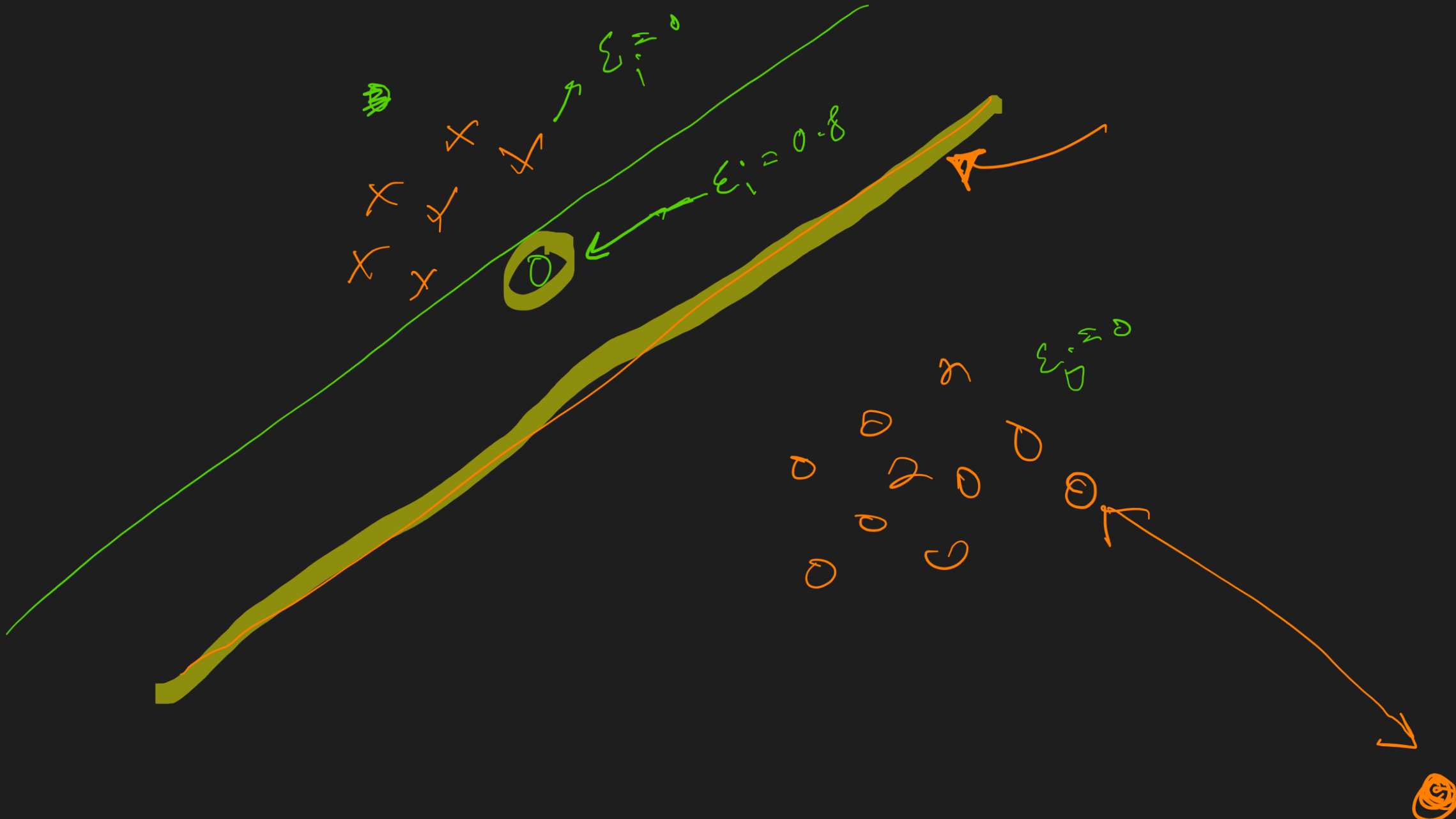
input: $(x_1, y_1), \dots, (x_m, y_m)$

solve:

$$(\underline{w}_0, \underline{b}_0) = \underset{(w, b)}{\operatorname{argmin}} \|w\|^2 \quad \text{s.t.} \quad \forall i, \quad y_i(\underbrace{\langle w, x_i \rangle}_{\delta} + \underbrace{b}_{\delta}) \geq 1 \quad (15.2)$$

$$\text{output: } \hat{w} = \frac{w_0}{\|w_0\|}, \quad \hat{b} = \frac{b_0}{\|w_0\|}$$

- SENSITIVE TO OUTLIERS...



SOFT-MARGIN SVM

Soft-SVM

input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

parameter: $\lambda > 0$

solve:

$$\min_{\mathbf{w}, b, \xi} \left(\lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right) \quad (15.4)$$

$$\text{s.t. } \forall i, \quad y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0$$

output: \mathbf{w}, b

EQUIVALENT FORM OF SOFT-MARGIN SVM

$$\rightarrow \min_{\mathbf{w}} \left(\lambda \|\mathbf{w}\|^2 + L_S^{\text{hinge}}(\mathbf{w}) \right),$$

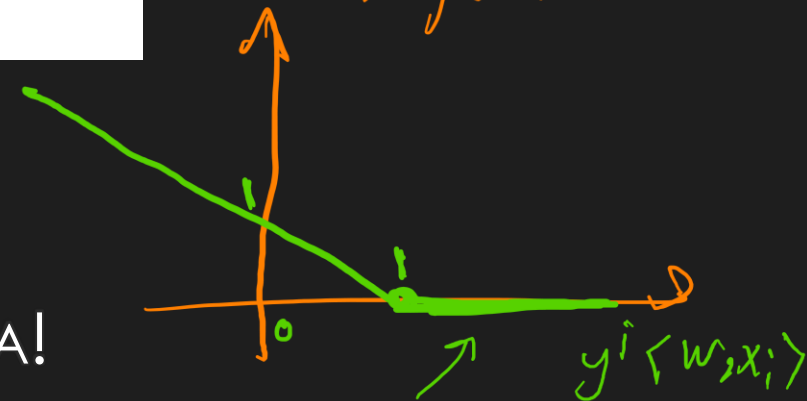
$$L_S^{\text{hinge}}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y \langle \mathbf{w}, \mathbf{x}_i \rangle\}.$$

Hinge loss

margin for data point i

Hinge loss

- NO CONSTRAINTS!
- REGULARIZATION!
 - SVMs ARE GOOD FOR HIGH-DIMENSIONAL DATA!



EXERCISE

- PROVE THAT THESE TWO FORMS OF SOFT-SVM ARE EQUIVALENT

$$\xi_i = \text{MAX}(0, 1 - y^i \langle \mathbf{w}, \mathbf{x}^i \rangle)$$

Soft-SVM

input: $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m)$

parameter: $\lambda > 0$

solve:

$$\min_{\mathbf{w}, b, \xi} \left(\lambda \|\mathbf{w}\|^2 + \frac{1}{m} \sum_{i=1}^m \xi_i \right) \quad (15.4)$$

$$\text{s.t. } \forall i, \quad y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0$$

output: \mathbf{w}, b

$$\min_{\mathbf{w}} \left(\lambda \|\mathbf{w}\|^2 + L_S^{\text{hinge}}(\mathbf{w}) \right),$$

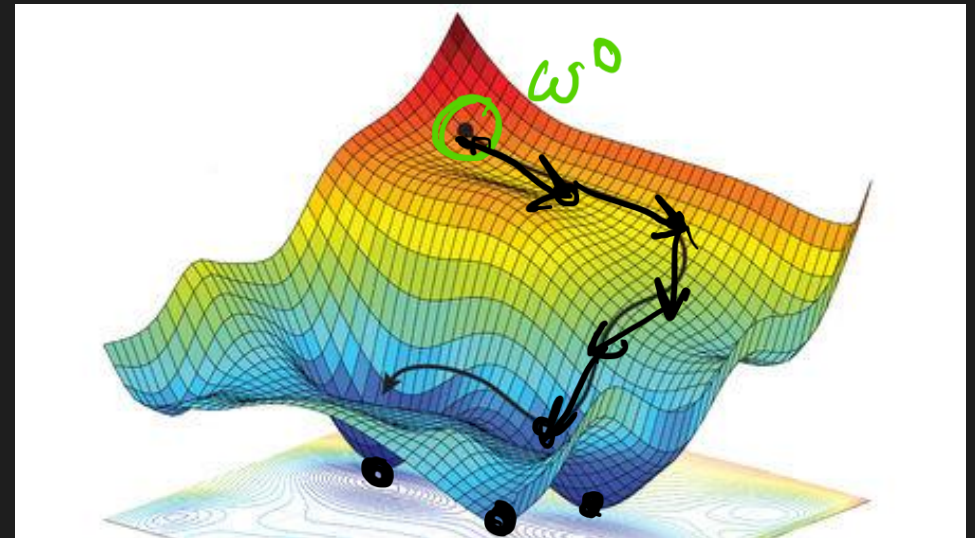
$$L_S^{\text{hinge}}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y \langle \mathbf{w}, \mathbf{x}_i \rangle\}.$$

OPTIMIZATION

- NO CONSTRAINTS!
- WE CAN USE SPECIAL-PURPOSE SVM SOLVERS
 - ...OR WE CAN JUST USE “GRADIENT DESCENT”!

$$\min_{\mathbf{w}} \left(\underbrace{\lambda \|\mathbf{w}\|^2}_{\text{regularization}} + \underbrace{L_S^{\text{hinge}}(\mathbf{w})}_{\text{loss}} \right),$$

$$L_S^{\text{hinge}}(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y \langle \mathbf{w}, \mathbf{x}_i \rangle\}.$$



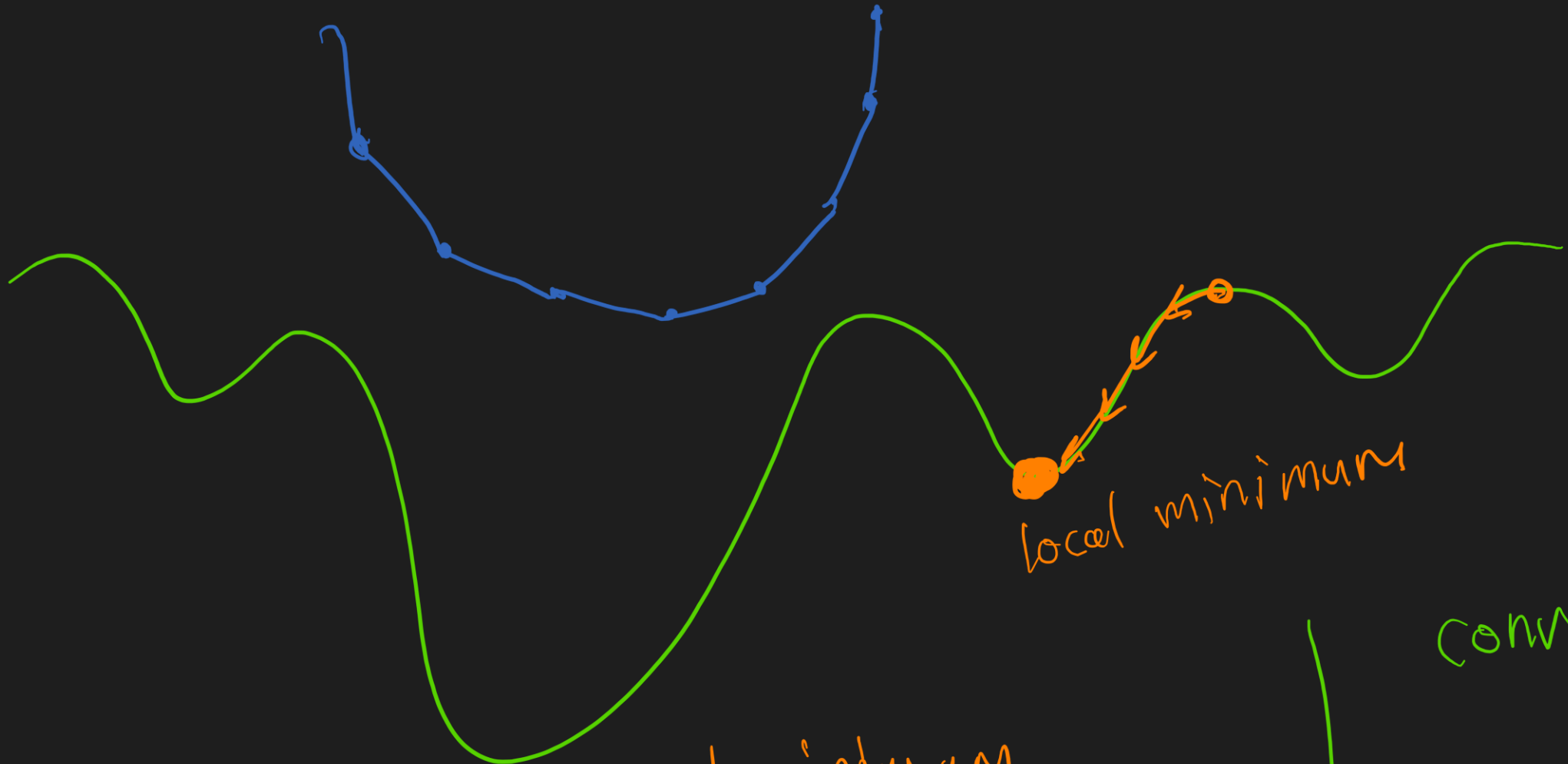
GRADIENT DESCENT

- MINIMIZE OBJECTIVE FUNCTION $E(w)$ WITH SINGLE REAL-VALUED OUTPUT?
- MOVE IN THE DIRECTION OF STEEPEST DESCENT
 - GRADIENT OF A FUNCTION AT A POINT IS THE DIRECTION OF STEEPEST ASCENT AT THAT POINT.

GRADIENT

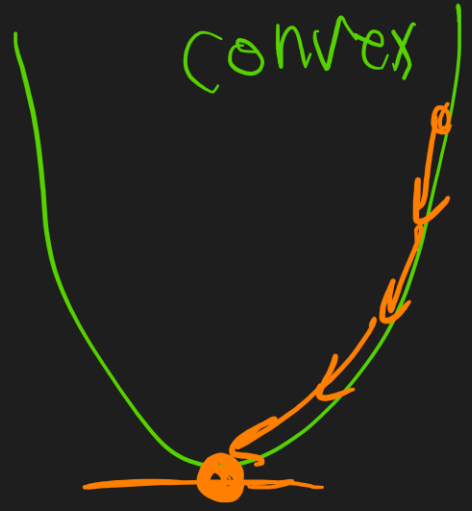
$$\nabla_w [E(w)] \approx \left[\frac{\partial E(w)}{\partial w_1}, \dots, \frac{\partial E(w)}{\partial w_D} \right]^T$$

- GRADIENT DESCENT
 - INITIALIZE w^0
 - FOR $t = 1, \dots$
 - $w^{t+1} = w^t - \alpha \nabla_w(E(w))$
 - TERMINATE AFTER A FEW ITERATIONS (OR UNTIL CONVERGENCE)
- α : STEP SIZE OR LEARNING RATE
 - AN IMPORTANT PARAMETER TO SET...CAN DEPEND ON t (α_t)



→ Global minimum

local minimum



convex

GRADIENT DESCENT: GUARANTEES

- CONVERGES TO A LOCAL MINIMUM/SADDLE POINT
 - WITH APPROPRIATE LEARNING RATE (CAN DEPEND ON t)
- NOT NECESSARILY GLOBAL MINIMUM
- GUARANTEED GLOBAL MINIMUM FOR CONVEX FUNCTIONS
 - SOFT SVM, LS, LOGISTIC REGRESSION (?)...
- MAY GET STUCK IN LOCAL MINIMUM FOR NON-CONVEX
 - WIDELY USED IN NEURAL NETWORKS (MORE ABOUT THIS LATER)

CALCULATING THE GRADIENT

- $E(w) = L(w) + \lambda \cdot \text{Reg}(w)$
 - E.G., FOR SOFT SVM WE HAVE $\text{Reg}(w) = \|w\|_2^2$
 - IT IS OFTEN EASY TO TAKE THE DERIVATIVE OF THE REGULARIZATION TERM
 - WHAT ABOUT THE LOSS TERM?
- $L(w) = \sum_i l(f_w(x^i), y^i)$
 - $\hat{y} = f_w(x)$ IS THE PREDICTED VALUE
 - $l(\hat{y}, y)$ IS THE LOSS FUNCTION
- $\nabla_w L(w) = \nabla_w \left(\sum_i l(f_w(x^i), y^i) \right) = \sum_i \nabla_w \left(l(f_w(x^i), y^i) \right)$
 - COMPUTE THE GRADIENT FOR EACH DATA POINT AND THEN SUM THEM UP
 - A BIT TOO SLOW.... MANY ITERATIONS

CALCULATING THE GRADIENT

- $\nabla_{\mathbf{w}}(L(\mathbf{w})) = \sum_i \nabla_{\mathbf{w}}(l(f_{\mathbf{w}}(x^i), y^i))$
- DIVIDE THE TRAINING DATA INTO A NUMBER OF “BATCHES”
- $S = S_1 \cup S_2 \cup \dots \cup S_m$
- $\nabla_{\mathbf{w}}^j = \sum_{(x,y) \in S_j} \nabla_{\mathbf{w}}(l(f_{\mathbf{w}}(x), y))$
- $\nabla_{\mathbf{w}}(E(\mathbf{w})) = \sum_j \nabla_{\mathbf{w}}^j$

MINI-BATCH GRADIENT DESCENT

- $\nabla_w(L(w)) = \sum_j \nabla_w^j$
- STANDARD GRADIENT DESCENT: $w^{t+1} = w^t - \alpha (\sum_j \nabla_w^j)$
- MINI-BATCH GRADIENT DESCENT:

- For $j = 1$ to m

- $w = w - \alpha (\nabla_w^j)$

- BENEFITS?
 - MEMORY
 - PARALLELIZATION
- POTENTIAL DISADVANTAGES
 - CAN BE LESS STABLE

Epoch: one time going over the whole dataset

STOCHASTIC GRADIENT DESCENT

Algorithm 8.1 Stochastic gradient descent (SGD) update

Require: Learning rate schedule $\epsilon_1, \epsilon_2, \dots$

Require: Initial parameter θ

$k \leftarrow 1$

while stopping criterion not met **do**

 Sample a minibatch of m examples from the training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ with corresponding targets $\mathbf{y}^{(i)}$.

 Compute gradient estimate: $\hat{\mathbf{g}} \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(\mathbf{x}^{(i)}; \theta), \mathbf{y}^{(i)})$

 Apply update: $\theta \leftarrow \theta - \epsilon_k \hat{\mathbf{g}}$

$k \leftarrow k + 1$

end while

- INSTEAD OF DIVIDING THE DATA INTO BATCHES, RANDOMLY SELECT A SUBSET AS A BATCH
- ANOTHER BENEFIT: ADDS A KIND OF RANDOMNESS

WHAT IS MISSING?

- $\nabla_{\mathbf{w}}(E(\mathbf{w})) = \sum_i \nabla_{\mathbf{w}} \left(l(f_{\mathbf{w}}(x^i), y^i) \right) + \nabla_{\mathbf{w}} \|\mathbf{w}\|_2^2$
 - HOW TO CALCULATE $\nabla_{\mathbf{w}} \left(l(f_{\mathbf{w}}(x^i), y^i) \right)$?
 - HOW TO CALCULATE $\nabla_{\mathbf{w}} \|\mathbf{w}\|_2^2$?
- WILL GET BACK TO THIS LATER BUT FOR NOW...
- “AUTOMATIC DIFFERENTIATION”

AUTOMATIC DIFFERENTIATION IN PYTORCH

```
# Define a linear SVM model
class LinearSVM(nn.Module):
    def __init__(self, input_size, num_classes):
        super(LinearSVM, self).__init__()
        self.fc = nn.Linear(input_size, num_classes, bias=True)

    def forward(self, x):
        # Flatten the image
        x = x.view(-1, input_size)
        return self.fc(x)
```

$$\hat{f}_w(x) = w^T x$$

```
# Initialize the model, loss function, and optimizer
model = LinearSVM(input_size, num_classes)
criterion = nn.MultiMarginLoss() # A Multi-class version of Hinge loss
optimizer = optim.SGD(model.parameters(), lr=learning_rate)
```

```
# Training the model
for epoch in range(num_epochs):
    for i, (images_batch, labels_batch) in enumerate(train_loader):
        optimizer.zero_grad() # Clear the gradients
        outputs = model(images_batch) # Forward pass
        loss = criterion(outputs, labels_batch) # Calculate loss
        loss.backward() # Backward pass
        optimizer.step() # Update weights
```