



Department of Computer Science
UNIVERSITY OF COLORADO **BOULDER**



Machine Learning: Chenhao Tan

University of Colorado Boulder

LECTURE 16

Slides adapted from Jordan Boyd-Graber, Chris Ketelsen

Logistics

- Homework 3 is due on Sunday!

Overview

A little bit of history

Linear classifiers and margin

Hard-margin SVM

- Formulation

- Theoretical Guarantees

Outline

A little bit of history

Linear classifiers and margin

Hard-margin SVM

Formulation

Theoretical Guarantees

Outline of CSCI 4622

We've already covered stuff in blue!

- Problem formulations: **classification**, **regression**
- Supervised techniques: **decision trees**, **nearest neighbors**, **perceptron**, **linear models**, **neural networks**, support vector machine, kernel methods
- Unsupervised techniques: clustering, linear dimensionality reduction
- “Meta-techniques”: ensembles, expectation-maximization
- Understanding ML: **limits of learning**, **practical issues**, bias & fairness
- Recurring themes: **(stochastic) gradient descent**

History lesson

- 1962: Rosenblatt, Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms
 - First neuron-based learning algorithm
 - “Could learning anything that you could program”

History lesson

- 1962: Rosenblatt, Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms
 - First neuron-based learning algorithm
 - “Could learning anything that you could program”
- 1969: Minsky & Papert, Perceptron: An Introduction to Computational Geometry
 - First real complexity analysis
 - Showed, in principle, many things that perceptrons can't learn to do
 - Shut down any interest in neural networks

History lesson

- 1986: Rumelhart, Hinton & Williams, Back Propagation
 - Overcame many difficulties raised by Minsky et al.
 - Neural networks wildly popular again (for a while)

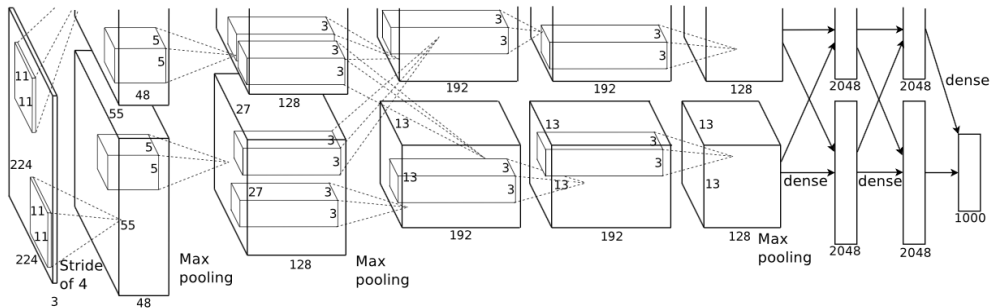
History lesson

- 1999-2005
 - Shift to Bayesian Methods
 - Best ideas from neural networks
 - Direct statistical computing
 - Support Vector Machines
 - Nice mathematical properties
 - Kernel tricks
 - A few people still playing with NNs
 - Bengio
 - Hinton
 - LeCun

History lesson

- 2005-2010
 - Core group continues to make improvements
 - Various tricks to make NNs practical
- 2010-present
 - BOOM!

AlexNet



Krizhevsky et al. [2012]

History lesson

Question: Why? What made neural networks amazing again?

- Massive datasets
- Computing power
- Algorithmic improvements

History lesson

Machine learning has a short history, but seems cyclic.
What is next?

Outline

A little bit of history

Linear classifiers and margin

Hard-margin SVM

Formulation

Theoretical Guarantees

Linear classifiers

We have already seen several:

- Logistic regression
- Perceptron

Definition: A linear classifier makes decisions by computing a linear combination of features of the form $\mathbf{w}^T \mathbf{x} + b$ and then classifies based on

$$\mathbf{w}^T \mathbf{x} + b \geq 0.$$

The decision boundary between the two classes is defined by $\mathbf{w}^T \mathbf{x} + b = 0$.

We estimate the weights and bias using the training data.

A linear classifier in 1D



- A linear classifier in 1D is a point x described by the equation $w_1 x_1 = -b$.

A linear classifier in 1D



- A linear classifier in 1D is a point x described by the equation $w_1 x_1 = -b$.
- $x_1 = -b/w_1$

A linear classifier in 1D



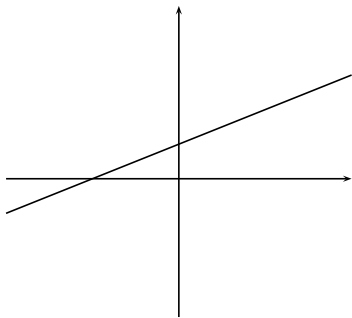
- A linear classifier in 1D is a point x described by the equation $w_1 x_1 = -b$.
- $x_1 = -b/w_1$
- Points (x_1) with $w_1 x_1 \geq -b$ are in the positive class.

A linear classifier in 1D



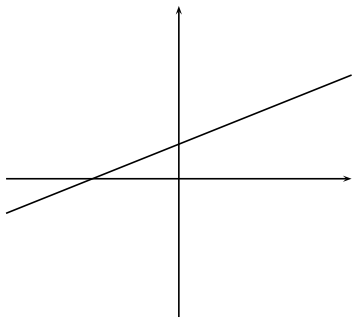
- A linear classifier in 1D is a point x described by the equation $w_1 x_1 = -b$.
- $x_1 = -b/w_1$
- Points (x_1) with $w_1 x_1 \geq -b$ are in the positive class.
- Points (x_1) with $w_1 x_1 < -b$ are in the negative class.

A linear classifier in 2D



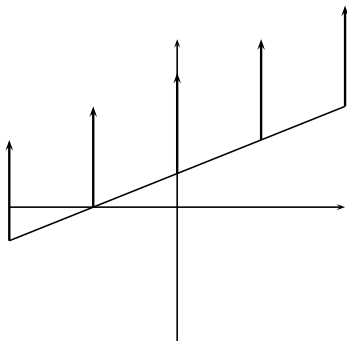
- A linear classifier in 2D is a line described by the equation $w_1x_1 + w_2x_2 = -b$.

A linear classifier in 2D



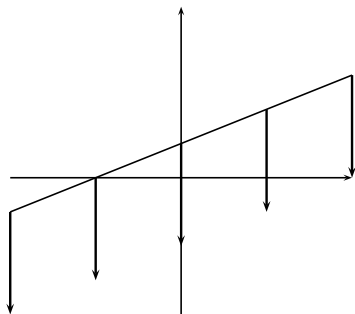
- A linear classifier in 2D is a line described by the equation $w_1x_1 + w_2x_2 = -b$.
- Example for a 2D linear classifier

A linear classifier in 2D



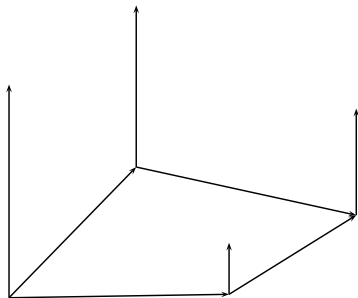
- A linear classifier in 2D is a line described by the equation $w_1x_1 + w_2x_2 = -b$.
- Example for a 2D linear classifier
- Points (x_1, x_2) with $w_1x_1 + w_2x_2 \geq -b$ are in the positive class.

A linear classifier in 2D



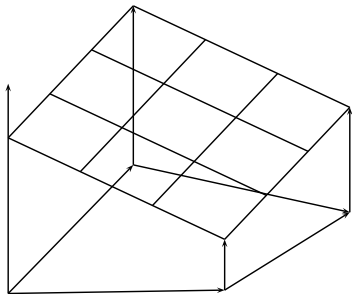
- A linear classifier in 2D is a line described by the equation $w_1x_1 + w_2x_2 = -b$.
- Example for a 2D linear classifier
- Points (x_1, x_2) with $w_1x_1 + w_2x_2 \geq -b$ are in the positive class.
- Points (x_1, x_2) with $w_1x_1 + w_2x_2 < -b$ are in the negative class.

A linear classifier in 3D



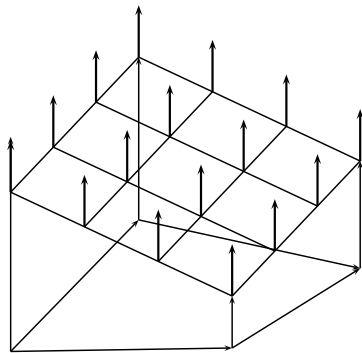
- A linear classifier in 3D is a plane described by the equation $w_1x_1 + w_2x_2 + w_3x_3 = -b$.

A linear classifier in 3D



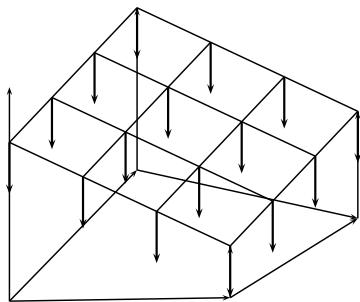
- A linear classifier in 3D is a plane described by the equation $w_1x_1 + w_2x_2 + w_3x_3 = -b$.
- Example for a 3D linear classifier

A linear classifier in 3D



- A linear classifier in 3D is a plane described by the equation $w_1x_1 + w_2x_2 + w_3x_3 = -b$.
- Example for a 3D linear classifier
- Points (x_1, x_2, x_3) with $w_1x_1 + w_2x_2 + w_3x_3 \geq -b$ are in the class c .

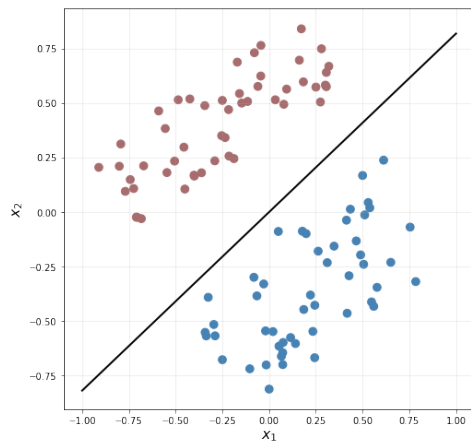
A linear classifier in 3D



- A linear classifier in 3D is a plane described by the equation $w_1x_1 + w_2x_2 + w_3x_3 = -b$.
- Example for a 3D linear classifier
- Points (x_1, x_2, x_3) with $w_1x_1 + w_2x_2 + w_3x_3 \geq -b$ are in the class c .
- Points (x_1, x_2, x_3) with $w_1x_1 + w_2x_2 + w_3x_3 < -b$ are in the complement class \bar{c} .

Pictorial definition of margin

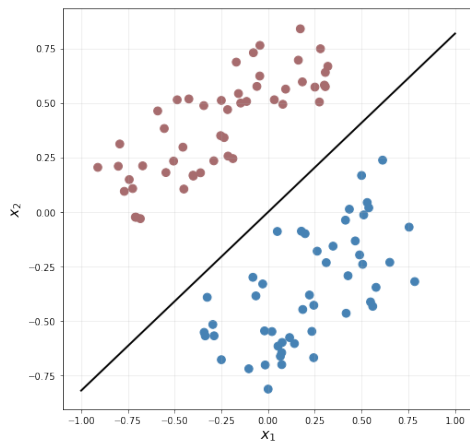
Assuming that the dataset is separable, margin is defined as the smallest distance from any data point to the decision boundary.



Pictorial definition of margin

Assuming that the dataset is separable, margin is defined as the smallest distance from any data point to the decision boundary.

What is the margin of the classifier on the right?



Outline

A little bit of history

Linear classifiers and margin

Hard-margin SVM

- Formulation

- Theoretical Guarantees

Which hyperplane?

- For linearly separable training sets: there are **infinitely** many separating hyperplanes.
- They all separate the training set perfectly . . .
- . . . but they behave differently on test data.
- Error rates on new data are low for some, high for others.
- How do we find a low-error separator?

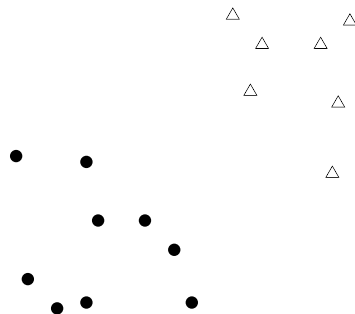
Support vector machines

SVMs: A kind of large-margin classifier

Find a decision boundary between two classes that is maximally far from any point in the training data (possibly discounting some points as outliers or noise).

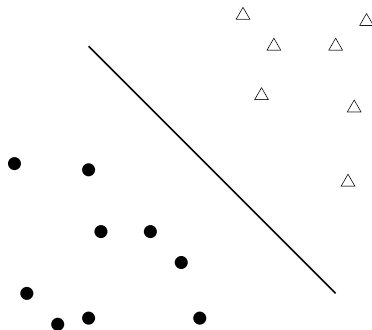
Support Vector Machines

- 2-class training data



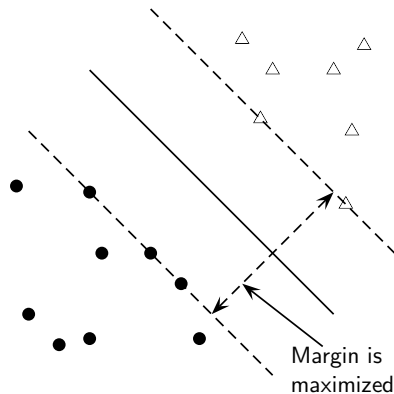
Support Vector Machines

- 2-class training data
- decision boundary \rightarrow **linear separator**



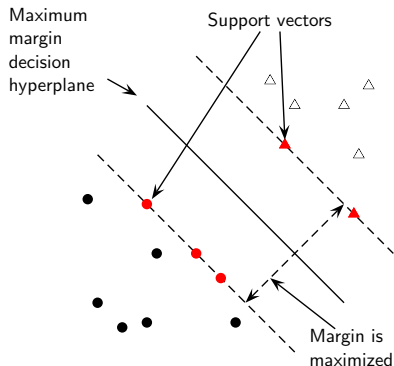
Support Vector Machines

- 2-class training data
- decision boundary \rightarrow **linear separator**
- criterion: being maximally far away from any data point \rightarrow determines classifier **margin**



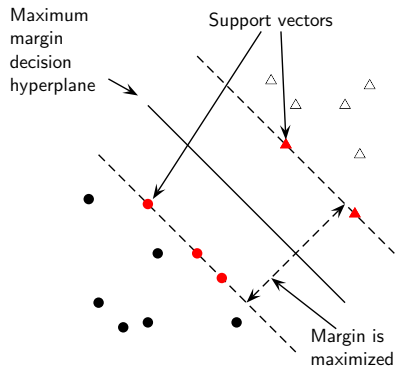
Support Vector Machines

- 2-class training data
- decision boundary → **linear separator**
- criterion: being maximally far away from any data point → determines classifier **margin**
- linear separator position defined by **support vectors**
- other points have no impact on the decision boundary



Why maximize the margin?

- Points near decision surface \rightarrow uncertain classification decisions
- A classifier with a large margin is always confident
- Gives classification safety margin (measurement or variation)
- Increased ability to correctly generalize to test data



Equation

- Equation of a hyperplane

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

- Distance of a point to hyperplane

$$\frac{|\mathbf{w} \cdot \mathbf{x} + b|}{\|\mathbf{w}\|}$$

- The margin ρ is given by

$$\rho \equiv \min_{(\mathbf{x}, y) \in S} \frac{|\mathbf{w} \cdot \mathbf{x} + b|}{\|\mathbf{w}\|} \equiv \frac{1}{\|\mathbf{w}\|}$$

Equation

- Equation of a hyperplane

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

- Distance of a point to hyperplane

$$\frac{|\mathbf{w} \cdot \mathbf{x} + b|}{\|\mathbf{w}\|}$$

- The margin ρ is given by

$$\rho \equiv \min_{(\mathbf{x}, y) \in S} \frac{|\mathbf{w} \cdot \mathbf{x} + b|}{\|\mathbf{w}\|} \equiv \frac{1}{\|\mathbf{w}\|}$$

This is because for any point on the marginal hyperplane, we can let $|\mathbf{w} \cdot \mathbf{x} + b| = 1$, and we would like to maximize the margin ρ .

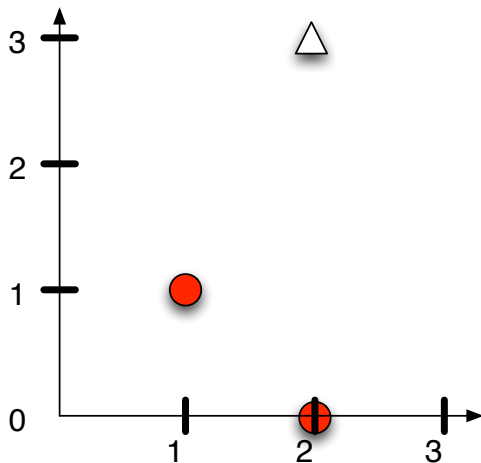
Optimization Problem

We want to find a weight vector \mathbf{w} and bias b that optimize

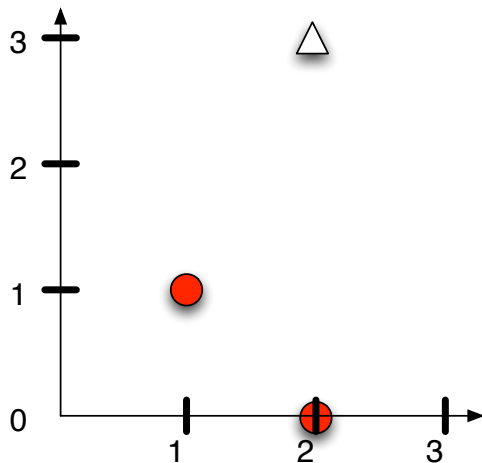
$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \forall i \in [1, n]$.

Find the maximum margin hyperplane



Find the maximum margin hyperplane



Which are the support vectors?

Walk through example: building an SVM over the data shown

Working geometrically:

- Set up system of equations

Walk through example: building an SVM over the data shown

Working geometrically:

- Set up system of equations

$$w_1 + w_2 + b = -1$$

$$\frac{3}{2}w_1 + 2w_2 + b = 0$$

$$2w_1 + 3w_2 + b = +1$$

Walk through example: building an SVM over the data shown

Working geometrically:

- Set up system of equations

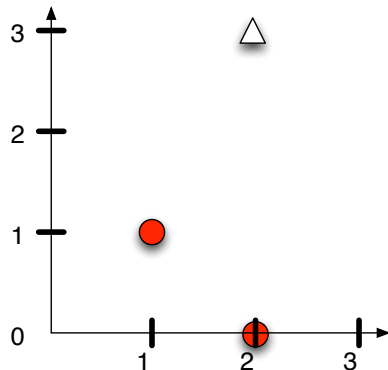
$$w_1 + w_2 + b = -1$$

$$\frac{3}{2}w_1 + 2w_2 + b = 0$$

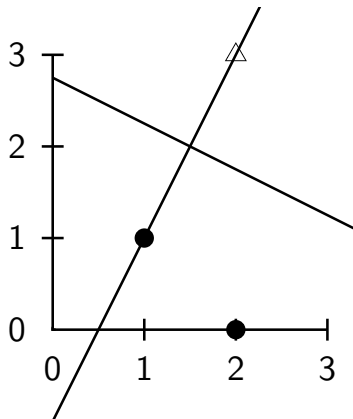
$$2w_1 + 3w_2 + b = +1$$

The SVM decision boundary is:

$$0 = \frac{2}{5}x + \frac{4}{5}y - \frac{11}{5}$$

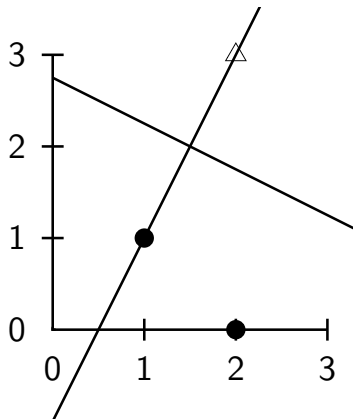


Canonical form



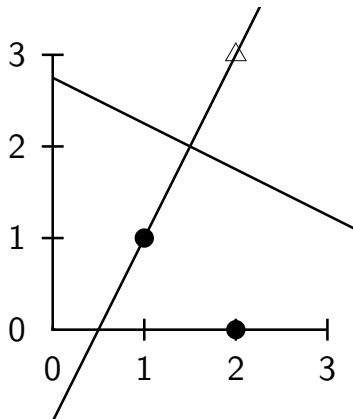
$$w_1x_1 + w_2x_2 + b = 0$$

Canonical form



$$.4x_1 + .8x_2 - 2.2 = 0$$

Canonical form



$$.4x_1 + .8x_2 - 2.2 = 0$$

- $.4 \cdot 1 + .8 \cdot 1 - 2.2 = -1$
- $.4 \cdot \frac{3}{2} + .8 \cdot 2 - 2.2 = 0$
- $.4 \cdot 2 + .8 \cdot 3 - 2.2 = +1$

What's the margin?

- Distance to closest point

What's the margin?

- Distance to closest point

$$\sqrt{\left(\frac{3}{2} - 1\right)^2 + (2 - 1)^2} = \frac{\sqrt{5}}{2}$$

What's the margin?

- Distance to closest point

$$\sqrt{\left(\frac{3}{2} - 1\right)^2 + (2 - 1)^2} = \frac{\sqrt{5}}{2}$$

- Margin computed from the weight vector

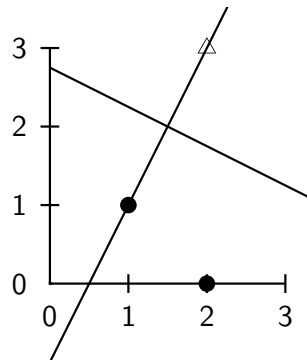
What's the margin?

- Distance to closest point

$$\sqrt{\left(\frac{3}{2} - 1\right)^2 + (2 - 1)^2} = \frac{\sqrt{5}}{2}$$

- Margin computed from the weight vector

$$\frac{1}{\|w\|} = \frac{1}{\sqrt{\left(\frac{2}{5}\right)^2 + \left(\frac{4}{5}\right)^2}} = \frac{1}{\sqrt{\frac{20}{25}}} = \frac{5}{\sqrt{5}\sqrt{4}} = \frac{\sqrt{5}}{2}$$



$$\frac{2}{5}x + \frac{4}{5}y - \frac{11}{5} = 0$$

Theoretical evidence that suggests SVMs will Work

- Leave-one-out error

Leave One Out Error (sketch)

Leave one out error is the error by using one point as your test set (averaged over all such points).

$$\hat{R}_{LOO} = \frac{1}{m} \sum_{i=1}^m \mathbb{1} [h_{s-\{x_i\}} \neq y_i] \quad (1)$$

Leave One Out Error (sketch)

Leave one out error is the error by using one point as your test set (averaged over all such points).

$$\hat{R}_{LOO} = \frac{1}{m} \sum_{i=1}^m \mathbb{1} [h_{s-\{x_i\}} \neq y_i] \quad (1)$$

This serves as an unbiased estimate of generalization error for samples of size $m - 1$:

Leave One Out Error (sketch)

Leave-one-out error is bounded by the number of support vectors.

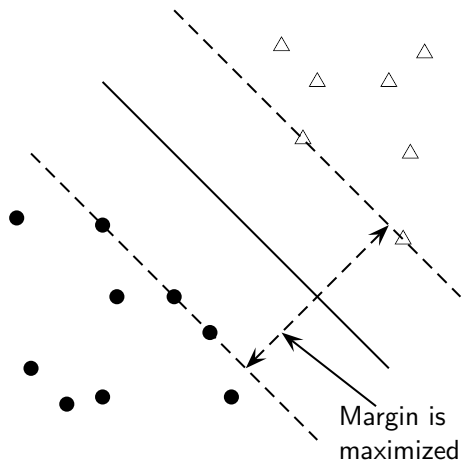
$$\mathbb{E}_{S \sim D^{m-1}} [R(h_S)] \leq \mathbb{E}_{S \sim D^m} \left[\frac{N_{SV}(S)}{m} \right] \quad (2)$$

Consider the held out error for x_i .

- If x_i was not a support vector, the answer doesn't change.
- If x_i was a support vector, it could change the answer; this is when we can have an error.

There are N_{SV} support vectors and thus N_{SV} possible errors.

Pictorial proof



References

Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems*, 2012.