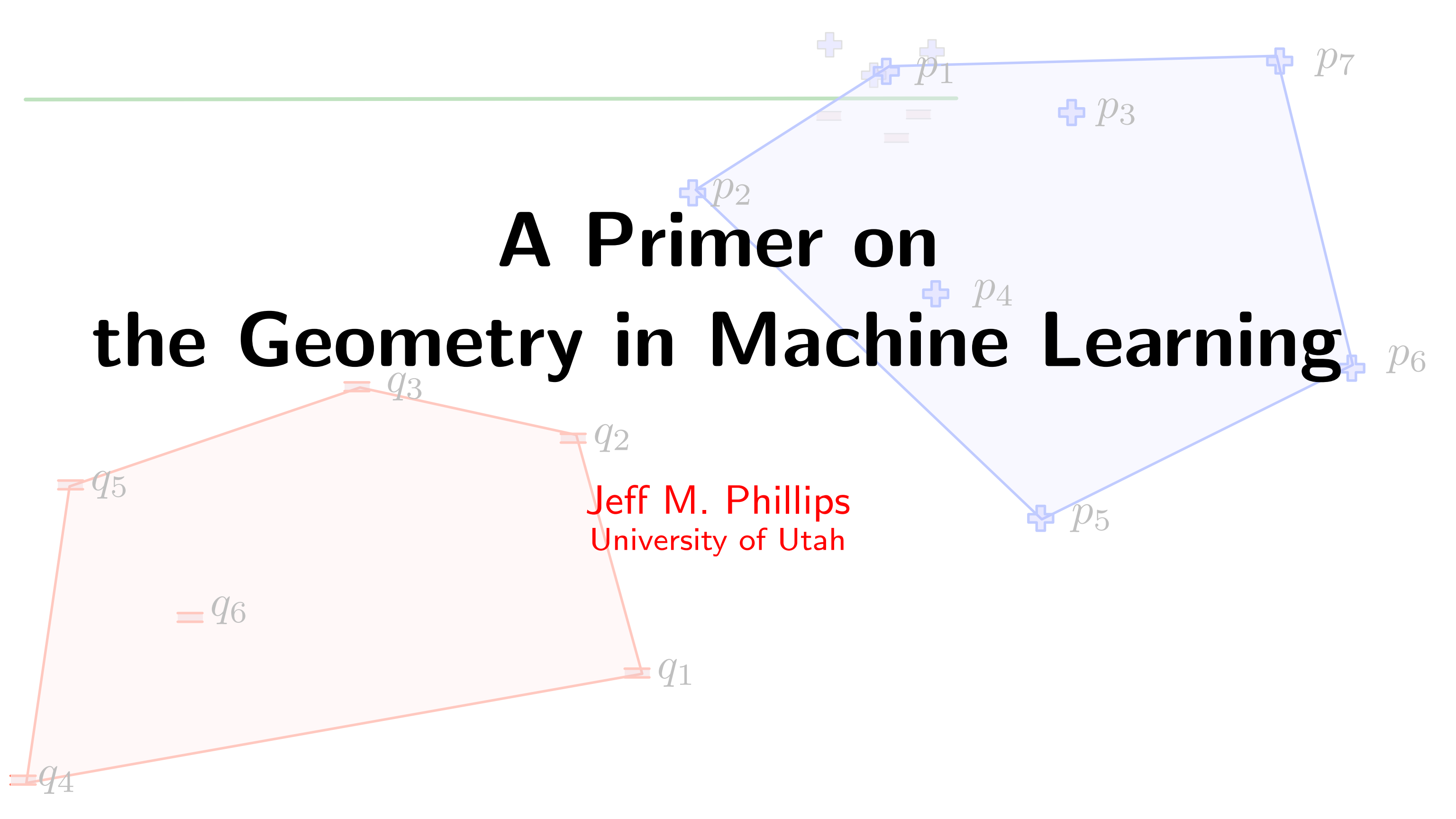# A Primer on
# the Geometry in Machine Learning

Jeff M. Phillips
University of Utah

$p_1$ $p_7$ $p_3$ $p_2$ $p_4$ $p_6$ $p_5$

$q_3$ $q_2$ $q_5$ $q_6$ $q_1$ $q_4$

# What is Machine Learning?

**supervised (has labels)**

**unsupervised (no labels)**

**class output**

**classification**

**clustering**

**value output**

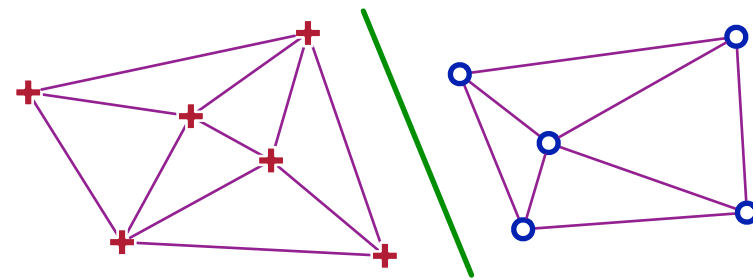**regression**
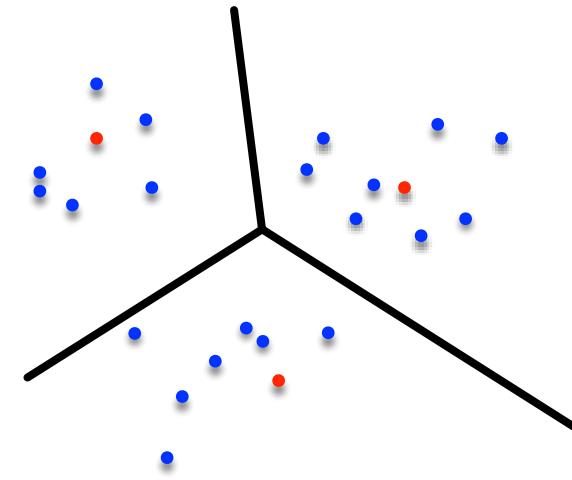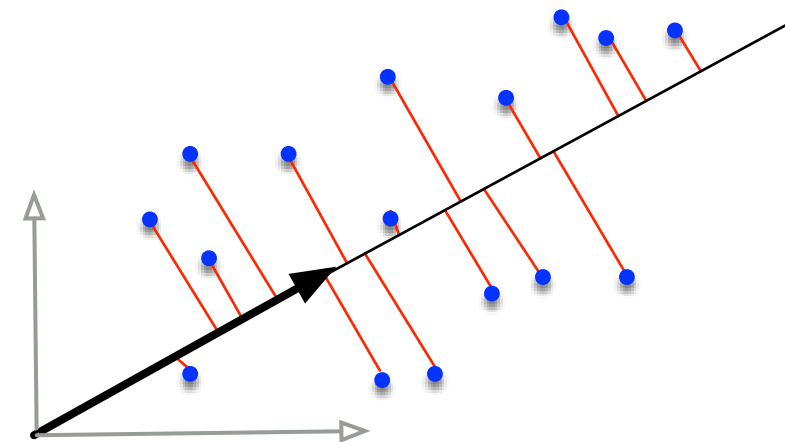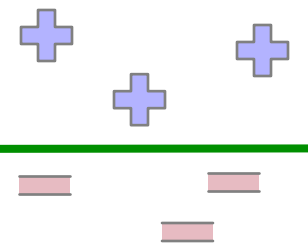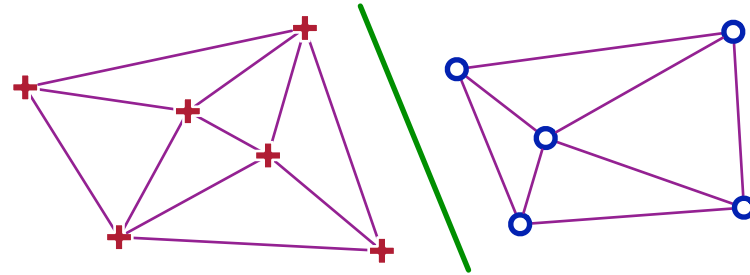
**dimensionality reduction**

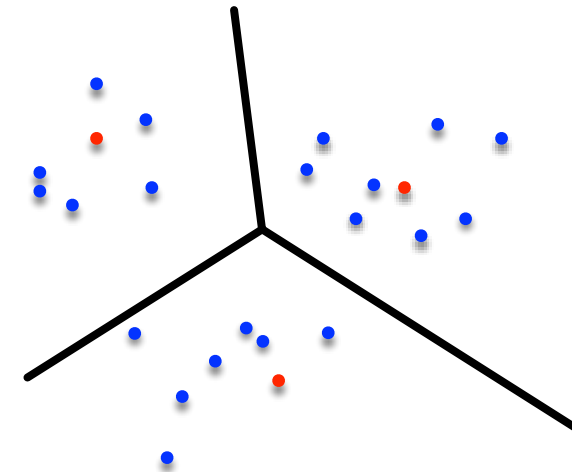# What is Machine Learning?

supervised (has labels)

unsupervised (no labels)

class
output

**classification**

**clustering**
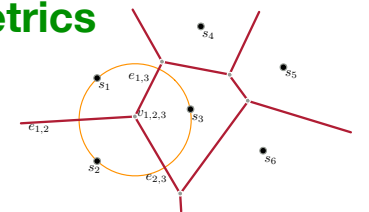
**Voronoi
metrics**

value
output

**regression**

**regularization**

**dimensionality reduction**

**linear projections**

# MATHEMATICAL FOUNDATIONS

## FOR

# DATA ANALYSIS

## JEFF M. PHILLIPS

http://www.cs.utah.edu/~jeffp/M4D/M4D.html

# What is Machine Learning?

Given $X \in \mathbb{R}^d$ with sign $\sigma : X \to \{-1, +1\}$.

Find separating halfspace $h \in \mathcal{H}$ so
$x \in h \Rightarrow \sigma(x) = +1$ and $x \notin h \Rightarrow \sigma(x) = -1$.

# What is Machine Learning?

Given $X \in \mathbb{R}^d$ with sign $\sigma : X \to \{-1, +1\}$.

Find separating halfspace $h \in \mathcal{H}$ so
$x \in h \Rightarrow \sigma(x) = +1$ and $x \notin h \Rightarrow \sigma(x) = -1$.

# What is Machine Learning?

Given $X \in \mathbb{R}^d$ with sign $\sigma : X \to \{-1, +1\}$.

Find separating halfspace $h \in \mathcal{H}$ so
$x \in h \Rightarrow \sigma(x) = +1$ and $x \notin h \Rightarrow \sigma(x) = -1$.

margin

# What is Machine Learning?

Given $X \in \mathbb{R}^d$ with sign $\sigma : X \to \{-1, +1\}$.

Find separating halfspace $h \in \mathcal{H}$ so
$x \in h \Rightarrow \sigma(x) = +1$ and $x \notin h \Rightarrow \sigma(x) = -1$.

**max margin**

# What is Machine Learning?

Given $X \in \mathbb{R}^d$ with sign $\sigma : X \to \{-1, +1\}$.

Find separating halfspace $h \in \mathcal{H}$ so
$x \in h \Rightarrow \sigma(x) = +1$ and $x \notin h \Rightarrow \sigma(x) = -1$.

max
margin

d+1
support
points

# What is Machine Learning?

Given $X \in \mathbb{R}^d$ with sign $\sigma : X \to \{-1, +1\}$.

Find separating halfspace $h \in \mathcal{H}$ so
$x \in h \Rightarrow \sigma(x) = +1$ and $x \notin h \Rightarrow \sigma(x) = -1$.

**How do we solve this problem?**

max
margin

d+1
support
points

# Polytope Distance

Define polytopes $P = \mathsf{CH}(x \in X \mid \sigma(x) = +1)$
$Q = \mathsf{CH}(x \in X \mid \sigma(x) = -1)$

$p_1$

$p_7$

$p_3$

$p_2$

$p_4$

$p_6$

$q_3$

$q_2$

$q_5$

$p_5$

$q_6$

$q_1$

$q_4$

# Polytope Distance

Define polytopes $P = \mathsf{CH}(x \in X \mid \sigma(x) = +1)$
$Q = \mathsf{CH}(x \in X \mid \sigma(x) = -1)$

# Polytope Distance

Define polytopes $P = \mathsf{CH}(x \in X \mid \sigma(x) = +1)$
$Q = \mathsf{CH}(x \in X \mid \sigma(x) = -1)$

Find $\arg\min_{p \in P, q \in Q} \|p - q\|$

# Polytope Distance

Define polytopes $P = \mathsf{CH}(x \in X \mid \sigma(x) = +1)$
$Q = \mathsf{CH}(x \in X \mid \sigma(x) = -1)$

Find $\arg\min_{p \in P, q \in Q} \|p - q\|$ ... a quadratic program in $\mathbb{R}^{2d}$

# Polytope Distance

Define polytopes $P = \mathsf{CH}(x \in X \mid \sigma(x) = +1)$
$Q = \mathsf{CH}(x \in X \mid \sigma(x) = -1)$

Find $\arg\min_{p \in P, q \in Q} \|p - q\|$ ... a quadratic program in $\mathbb{R}^{2d}$

Assume all $x \in X$ have $\|x\| \leq 1$
Let $\gamma = \min_{p \in P, q \in Q} \|p - q\|$

# Polytope Distance

Define polytopes $P = \mathsf{CH}(x \in X \mid \sigma(x) = +1)$
$Q = \mathsf{CH}(x \in X \mid \sigma(x) = -1)$

Find $\arg\min_{p \in P, q \in Q} \|p - q\|$ ... a quadratic program in $\mathbb{R}^{2d}$

Assume all $x \in X$ have $\|x\| \leq 1$
Let $\gamma = \min_{p \in P, q \in Q} \|p - q\|$

Iterate $1/(\varepsilon\gamma^2)$ steps,
find $\hat{p} \in P$ and $\hat{q} \in Q$ so
$(1 - \varepsilon)\|\hat{p} - \hat{q}\| \leq \gamma$

# Polytope Distance



Define polytopes $P = \mathsf{CH}(x \in X \mid \sigma(x) = +1)$
$Q = \mathsf{CH}(x \in X \mid \sigma(x) = -1)$

Find $\arg\min_{p \in P, q \in Q} \|p - q\|$ ... a quadratic program in $\mathbb{R}^{2d}$

Assume all $x \in X$ have $\|x\| \leq 1$
Let $\gamma = \min_{p \in P, q \in Q} \|p - q\|$

Iterate $1/(\varepsilon\gamma^2)$ steps,
find $\hat{p} \in P$ and $\hat{q} \in Q$ so
$(1 - \varepsilon)\|\hat{p} - \hat{q}\| \leq \gamma$

# Polytope Distance

Define polytopes $P = \mathsf{CH}(x \in X \mid \sigma(x) = +1)$
$Q = \mathsf{CH}(x \in X \mid \sigma(x) = -1)$

Find $\arg\min_{p \in P, q \in Q} \|p - q\|$ ... a quadratic program in $\mathbb{R}^{2d}$

Assume all $x \in X$ have $\|x\| \leq 1$
Let $\gamma = \min_{p \in P, q \in Q} \|p - q\|$

Iterate $1/(\varepsilon \gamma^2)$ steps,
find $\hat{p} \in P$ and $\hat{q} \in Q$ so
  $(1 - \varepsilon)\|\hat{p} - \hat{q}\| \leq \gamma$

# Polytope Distance

Define polytopes $P = \mathsf{CH}(x \in X \mid \sigma(x) = +1)$
$Q = \mathsf{CH}(x \in X \mid \sigma(x) = -1)$

Find $\arg\min_{p \in P, q \in Q} \|p - q\|$ ... a quadratic program in $\mathbb{R}^{2d}$

Assume all $x \in X$ have $\|x\| \leq 1$
Let $\gamma = \min_{p \in P, q \in Q} \|p - q\|$

Iterate $1/(\varepsilon \gamma^2)$ steps,
find $\hat{p} \in P$ and $\hat{q} \in Q$ so
$(1 - \varepsilon)\|\hat{p} - \hat{q}\| \leq \gamma$

# Polytope Distance

Define polytopes $P = \mathsf{CH}(x \in X \mid \sigma(x) = +1)$
$Q = \mathsf{CH}(x \in X \mid \sigma(x) = -1)$

Find $\arg\min_{p \in P, q \in Q} \|p - q\|$ ... a quadratic program in $\mathbb{R}^{2d}$

Assume all $x \in X$ have $\|x\| \leq 1$
Let $\gamma = \min_{p \in P, q \in Q} \|p - q\|$

Iterate $1/(\varepsilon\gamma^2)$ steps,
find $\hat{p} \in P$ and $\hat{q} \in Q$ so
$(1 - \varepsilon)\|\hat{p} - \hat{q}\| \leq \gamma$

$p_1$
$p_7$
$p_3$
$p_2$
$p_4$
$p_6$
$p_5$

$q_3$
$q_2$
$q_5$
$q_6$
$q_1$
$q_4$

**Gilbert 1966**
**Clarkson 2008**
**Gartner+Jaggi 2009**

# Force $h$ through Origin

1: Map $p \in \mathbb{R}^d$ to $p' \in \mathbb{R}^{d+1}$ as $p' = (p_1, p_2, \ldots, p_d, 1)$.

$p_1$

$p_7$

$p_3$

$p_2$

$p_4$

$p_6$

$q_3$

$q_2$

$q_5$

$p_5$

$q_6$

$q_1$

$q_4$

# Force $h$ through Origin

1: Map $p \in \mathbb{R}^d$ to $p' \in \mathbb{R}^{d+1}$ as $p' = (p_1, p_2, \ldots, p_d, 1)$.

# Force $h$ through Origin

1: Map $p \in \mathbb{R}^d$ to $p' \in \mathbb{R}^{d+1}$ as $p' = (p_1, p_2, \ldots, p_d, 1)$.

2: Force halfspace boundary $\partial h$ to include origin.

# Perceptron Algorithm

Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \leq 1$
(2) halfspace $h \in \mathcal{H}_0$ goes through origin

# Perceptron Algorithm

Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \leq 1$
(2) halfspace $h \in \mathcal{H}_0$ goes through origin

Algorithm: **(Rosenblatt 1958)**
    choose $w = \sigma(x)x$
    **for** $i = 1$ **to** $1/\gamma^2$ steps
    $x' = $ any $x \in X$ s.t. $\langle \sigma(x)x, w \rangle < 0$
    $w = w + \sigma(x')x'$

# Perceptron Algorithm

Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \leq 1$
(2) halfspace $h \in \mathcal{H}_0$ goes through origin

Algorithm:
    choose $w = \sigma(x)x$
    **for** $i = 1$ **to** $1/\gamma^2$ steps
        $x' = $ any $x \in X$ s.t. $\langle \sigma(x)x, w \rangle < 0$
        $w = w + \sigma(x')x'$

# Perceptron Algorithm

Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \leq 1$
(2) halfspace $h \in \mathcal{H}_0$ goes through origin

Algorithm:
    choose $w = \sigma(x)x$
    **for** $i = 1$ **to** $1/\gamma^2$ steps
        $x' = $ any $x \in X$ s.t. $\langle \sigma(x)x, w \rangle < 0$
        $w = w + \sigma(x')x'$

# Perceptron Algorithm

Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \leq 1$

(2) halfspace $h \in \mathcal{H}_0$ goes through origin

Algorithm:

    choose $w = \sigma(x)x$

    **for** $i = 1$ **to** $1/\gamma^2$ steps

        $x' = $ any $x \in X$ s.t. $\langle \sigma(x)x, w \rangle < 0$

        $w = w + \sigma(x')x'$

# Perceptron Algorithm

Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \leq 1$
(2) halfspace $h \in \mathcal{H}_0$ goes through origin

Algorithm:
    choose $w = \sigma(x)x$
    **for** $i = 1$ **to** $1/\gamma^2$ steps
      $x' =$ any $x \in X$ s.t. $\langle \sigma(x)x, w \rangle < 0$
      $w = w + \sigma(x')x'$

# Perceptron Algorithm

Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \leq 1$
(2) halfspace $h \in \mathcal{H}_0$ goes through origin

Algorithm:
  choose $w = \sigma(x)x$
  **for** $i = 1$ **to** $1/\gamma^2$ steps
    $x' = $ any $x \in X$ s.t. $\langle \sigma(x)x, w \rangle < 0$
    $w = w + \sigma(x')x'$

# Perceptron Algorithm

Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \le 1$
(2) halfspace $h \in \mathcal{H}_0$ goes through origin

Algorithm:
  choose $w = \sigma(x)x$
  **for** $i = 1$ **to** $1/\gamma^2$ steps
    $x' =$ any $x \in X$ s.t. $\langle \sigma(x)x, w \rangle < 0$
    $w = w + \sigma(x')x'$

# Perceptron Algorithm



Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \leq 1$
(2) halfspace $h \in \mathcal{H}_0$ goes through origin

Algorithm:
    choose $w = \sigma(x)x$
    **for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps
      $x' = \arg\min \langle \sigma(x)x, w \rangle$
      $w = w + \sigma(x')x'$



$w$

# Perceptron Algorithm

Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \le 1$
(2) halfspace $h \in \mathcal{H}_0$ goes through origin

Algorithm:
    choose $w = \sigma(x)x$
    **for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps
        $x' = \arg\min \langle \sigma(x)x, w \rangle$
        $w = w + \sigma(x')x'$

**What to do if the polytopes intersect?**
**it is ``non-separable'' ?**

# Perceptron Algorithm

Assume (1) $x \in X \subset \mathbb{R}^d$ has $\|x\| \leq 1$
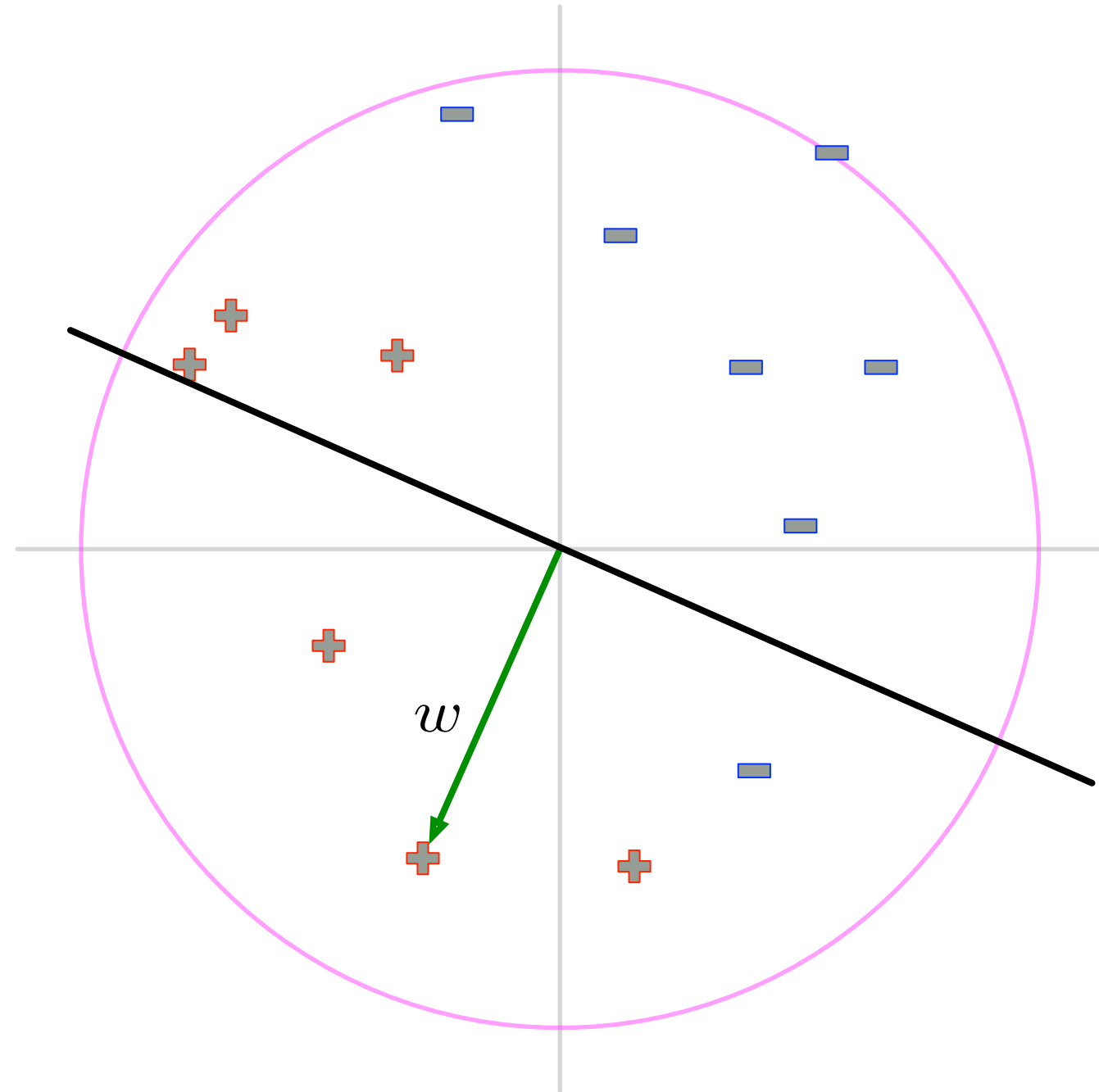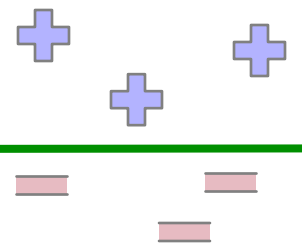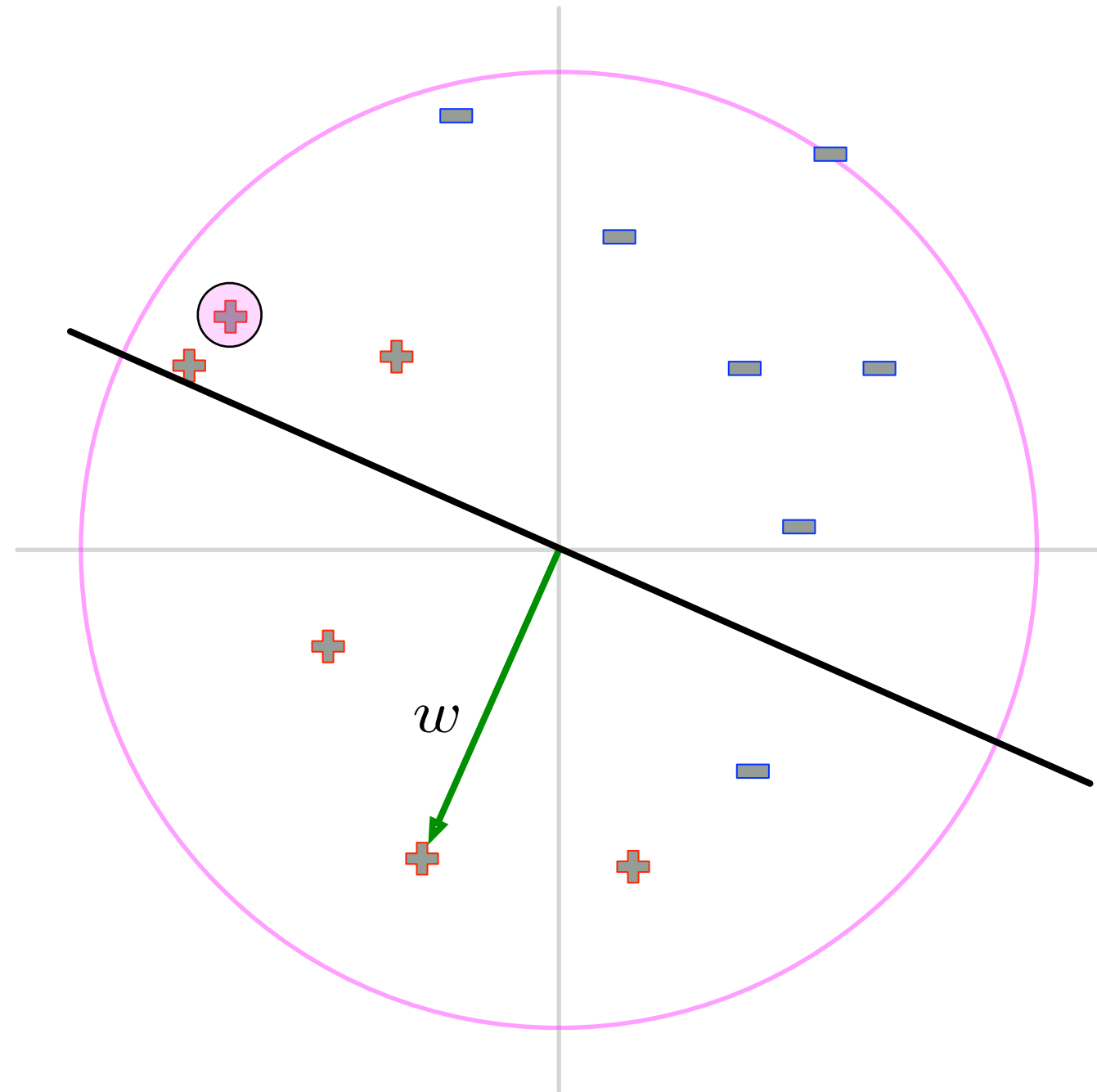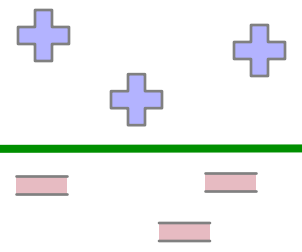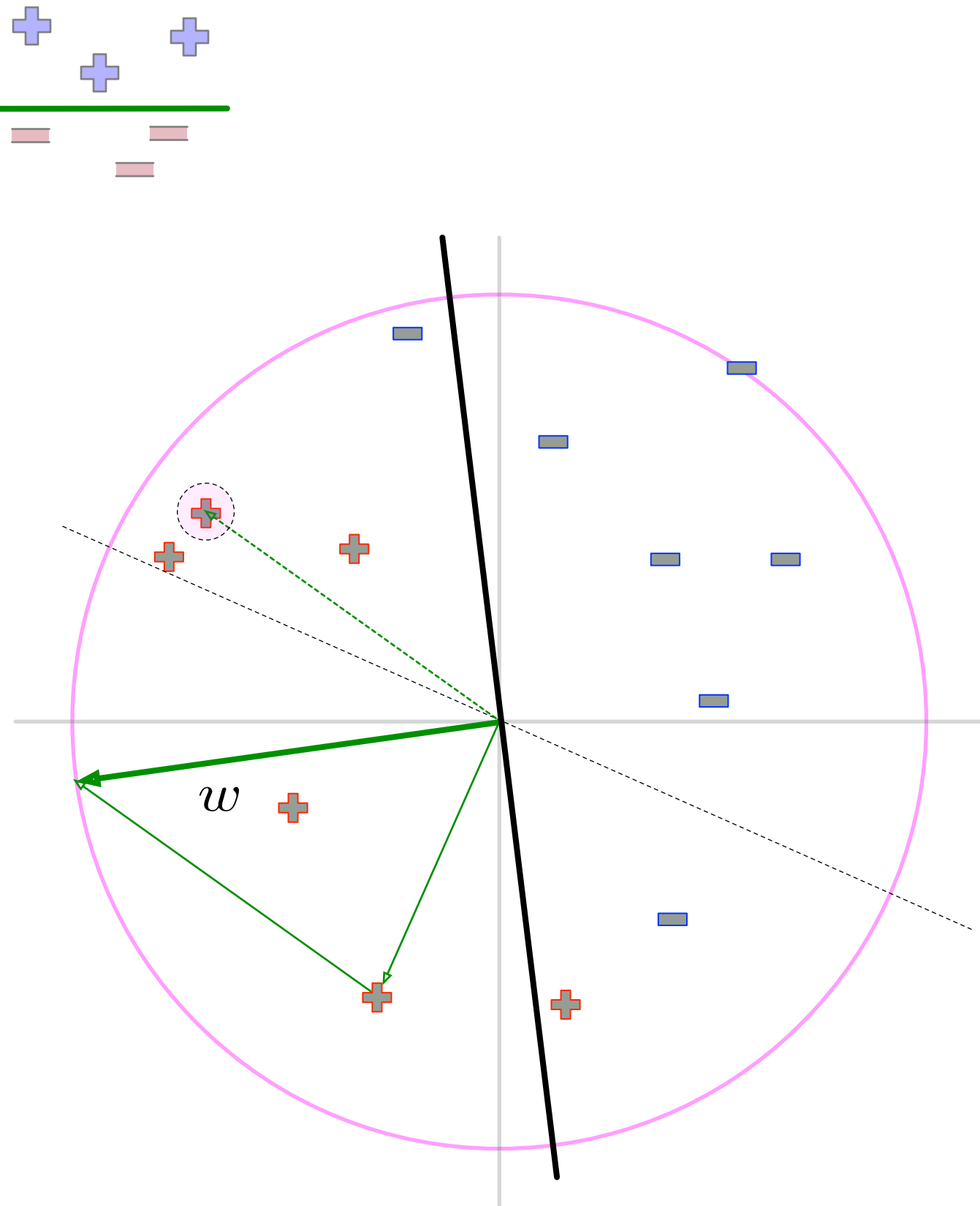(2) halfspace $h \in \mathcal{H}_0$ goes through origin

Algorithm:
    choose $w = \sigma(x)x$
    **for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps
      $x' = \arg\min \langle \sigma(x)x, w \rangle$
      $w = w + \sigma(x')x'$

**What to do if the polytopes intersect?**
    **it is "non-separable" ?**

    **(1) kernels**
    **(2) penalize it**

# Kernel Perceptron Algorithm

Algorithm: (Kernel SVM)

    choose $w = \sigma(x)x$

    **for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps

        $x' = \arg\min \sigma(x) \langle x, w \rangle$

        $w = w + \sigma(x')x'$

# Kernel Perceptron Algorithm

Algorithm: (Kernel SVM)

choose $w = \sigma(x)x$

**for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps

$x' = \arg\min \sigma(x)K(x,w)$

$w = w \oplus \sigma(x')x'$

$K(x,p) = \exp(-\|x-p\|^2)$

$K(x,p) = \exp(-\|x-p\|)$

$K(x,p) = (1 + \langle x,p\rangle)^r$

$p_1$        $p_7$

$p_3$

$p_2$

$p_4$

$q_3$        $p_6$

$q_2$

$q_5$

$p_5$

$q_6$

$q_1$

$q_4$

# Kernel Perceptron Algorithm
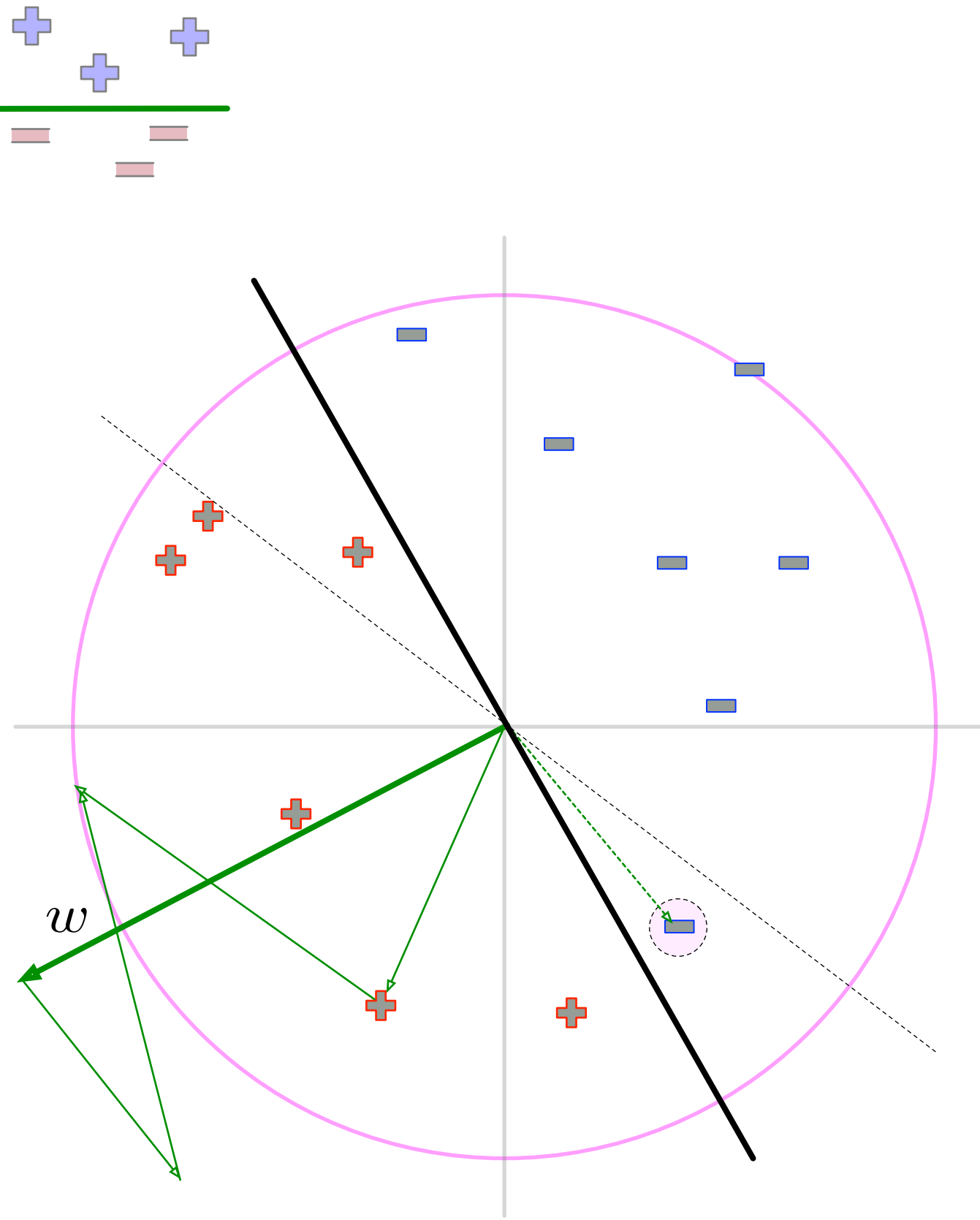
Algorithm: (Kernel SVM)

    choose $w = \sigma(x)x$

    **for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps

        $x' = \arg\min \sigma(x)\mathrm{KDE}_W(x)$

        $\mathrm{KDE}_W = \mathrm{KDE}_W + \sigma(x')K(x', \cdot)$

$K(x, p) = \exp(-\|x - p\|^2)$
$K(x, p) = \exp(-\|x - p\|)$
$K(x, p) = (1 + \langle x, p \rangle)^r$

kernel density estimate:
$\mathrm{KDE}_W(x) = \frac{1}{|W|} \sum_{w \in W} K(w, x)$

$p_1$      $p_7$

$p_3$

$p_2$

$p_4$

$q_3$

$p_6$

$q_2$

$q_5$

$p_5$

$q_6$

$q_1$

$q_4$

# Kernel Perceptron Algorithm

Algorithm: (Kernel SVM)
$\quad$ choose $w = \sigma(x)x$
$\quad$ **for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps
$\quad\quad x' = \arg\min \sigma(x)\mathrm{KDE}_W(x)$
$\quad\quad \mathrm{KDE}_W = \mathrm{KDE}_W + \sigma(x')K(x', \cdot)$

$K(x, p) = \exp(-\|x - p\|^2)$
$K(x, p) = \exp(-\|x - p\|)$
$K(x, p) = (1 + \langle x, p \rangle)^r$

kernel density estimate:
$\mathrm{KDE}_W(x) = \frac{1}{|W|} \sum_{w \in W} K(w, x)$

# Sketched Kernel SVM Algorithm

Algorithm:

Define mapping $\phi : X \to \mathbb{R}^\rho$

choose $w = \sigma(x)\phi(x)$

**for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps

$x' = \arg\min \sigma(x)\langle \phi(x), w \rangle$

$w = w + \sigma(x')\phi(x')$

# Sketched Kernel SVM Algorithm

Algorithm:

Define mapping $\phi : X \to \mathbb{R}^\rho$

    choose $w = \sigma(x)\phi(x)$

    **for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps
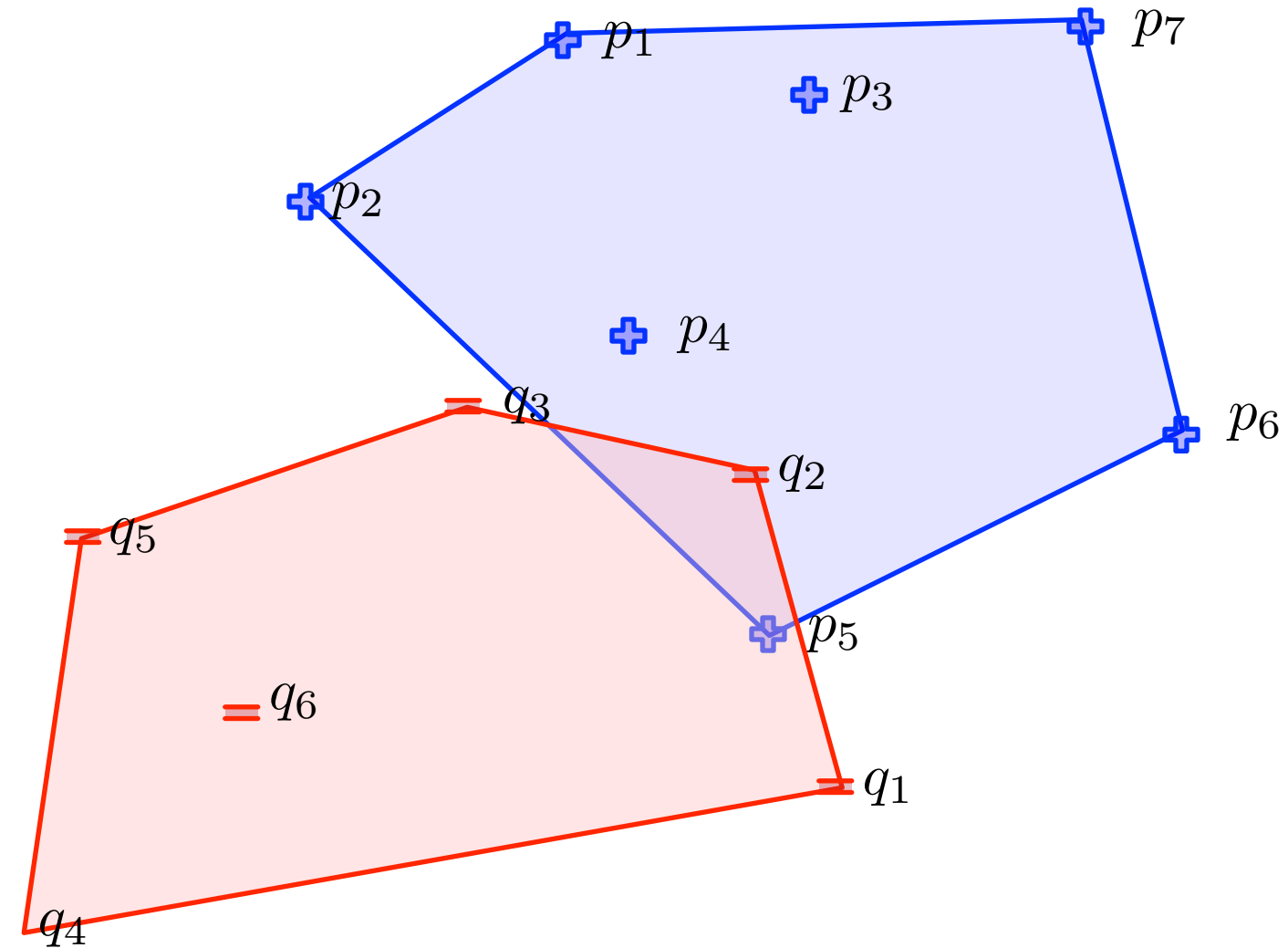
        $x' = \arg\min \sigma(x)\langle\phi(x), w\rangle$

        $w = w + \sigma(x')\phi(x')$

$K(x, p) = (1 + \langle x, p \rangle)^r \in \mathbb{R}^{O(rd)}$

$p_1$      $p_7$

$p_3$

$p_2$

$p_4$

$q_5$    $q_3$      $p_6$

$q_2$

$p_5$

$q_6$

$q_1$

$q_4$

# Sketched Kernel SVM Algorithm

Algorithm:

<mark>Define mapping $\phi : X \to \mathbb{R}^\rho$</mark>

    choose $w = \sigma(x)\phi(x)$

    **for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps

        $x' = \arg\min \sigma(x)\langle \phi(x), w \rangle$

        $w = w + \sigma(x')\phi(x')$

$K(x, p) = (1 + \langle x, p \rangle)^r \in \mathbb{R}^{O(rd)}$

$K(x, p) = \exp(-\|x - p\|^2)$ maps approximately to $\mathbb{R}^\rho$

e.g., Rahimi+Recht [NeurIPS07] additive $\varepsilon$ in $\rho = \tilde{O}(1/\varepsilon^2)$.

# Sketched Kernel SVM Algorithm

Algorithm:

Define mapping $\phi : X \to \mathbb{R}^\rho$

    choose $w = \sigma(x)\phi(x)$

    **for** $i = 1$ **to** $1/\varepsilon\gamma^2$ steps

      $x' = \arg\min \sigma(x)\langle \phi(x), w \rangle$
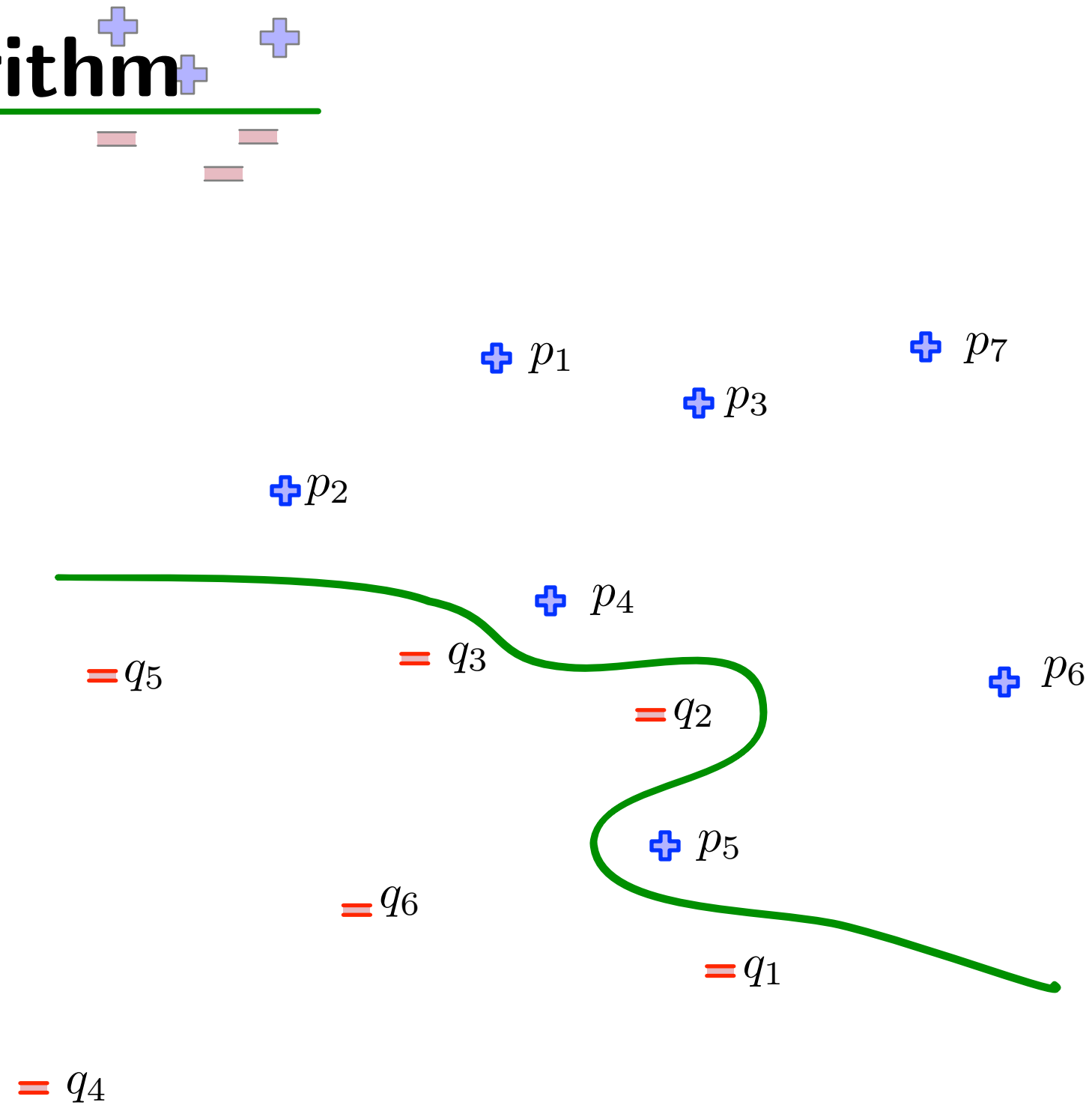
      $w = w + \sigma(x')\phi(x')$

$K(x, p) = (1 + \langle x, p \rangle)^r \in \mathbb{R}^{O(rd)}$

$K(x, p) = \exp(-\|x - p\|^2)$ maps approximately to $\mathbb{R}^\rho$

e.g., Rahimi+Recht [NeurIPS07] additive $\varepsilon$ in $\rho = \tilde{O}(1/\varepsilon^2)$.

$\phi(x) = K(x, \cdot)$ is a element of function space RKHS $\mathcal{H}_K$.

kernel density estimate:

$$\mathrm{KDE}_W = \frac{1}{|W|} \sum_{w \in W} \phi(w).$$

# Loss Functions

Set up penalty for misclassified points

$$f_X(w) = \sum_{x \in X} \ell(w, x) + \text{prior}(w)$$

loss $\ell_i = \ell(w, x_i) = \ell(z_i)$
with $z_i = \sigma(x_i)\langle w, x_i \rangle$

# Loss Functions

Set up penalty for misclassified points

$$f_X(w) = \sum_{x \in X} \ell(w, x) + \mathsf{prior}(w)$$

loss $\ell_i = \ell(w, x_i) = \ell(z_i)$
with $z_i = \sigma(x_i)\langle w, x_i \rangle$

# Loss Functions

Set up penalty for misclassified points

$$f_X(w) = \sum_{x \in X} \ell(w, x) + \mathsf{prior}(w)$$

loss $\ell_i = \ell(w, x_i) = \ell(z_i)$
with $z_i = \sigma(x_i)\langle w, x_i \rangle$

$p_1$  $p_7$

$p_3$

$p_2$

$p_4$

$q_3$

$q_2$

$q_5$

$p_6$

$p_5$

$q_6$

$q_4$

$q_1$

$\Delta$

hinge

$\ell_i$

smoothed hinge

$\ell_i$

squared hinge
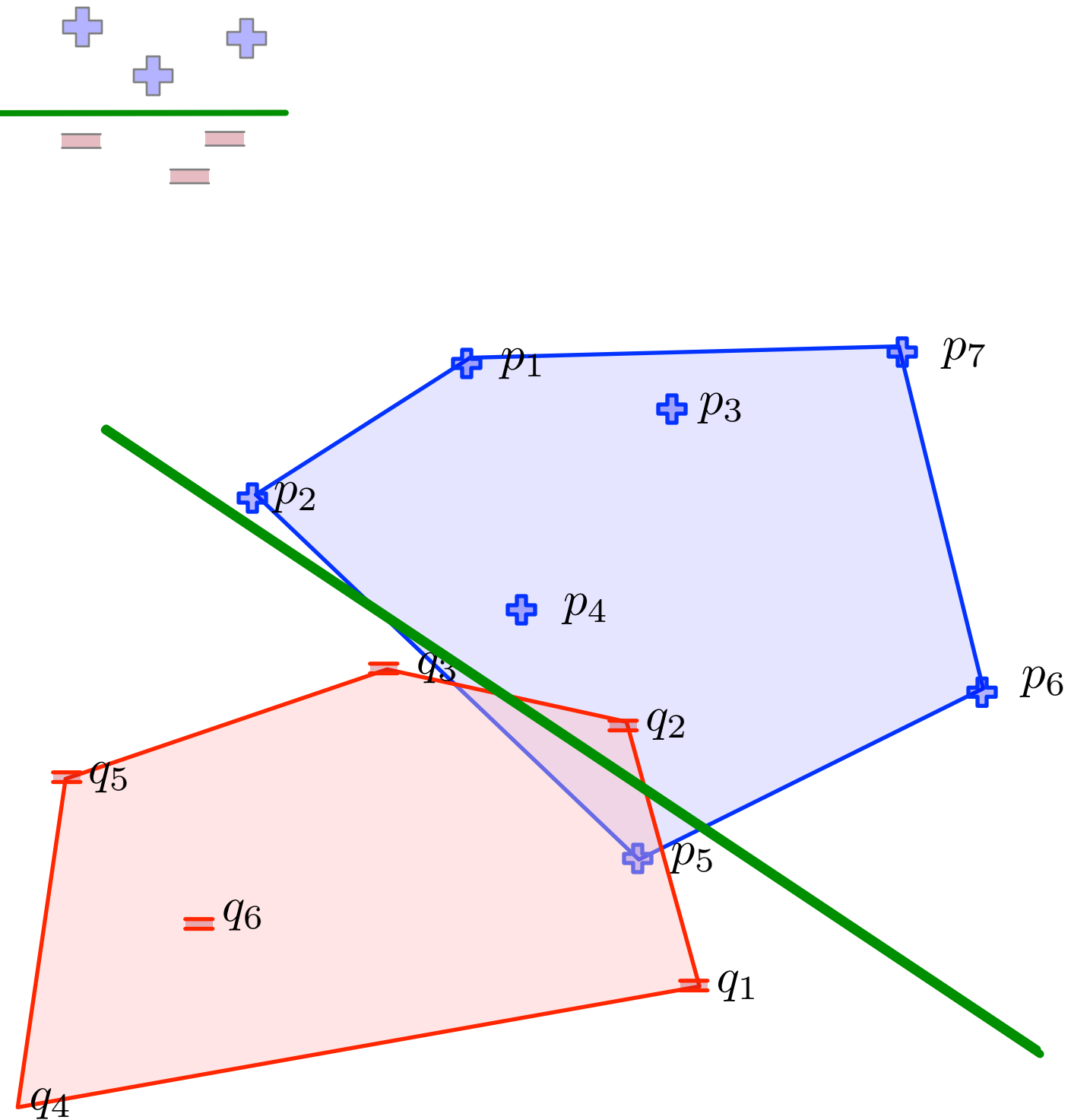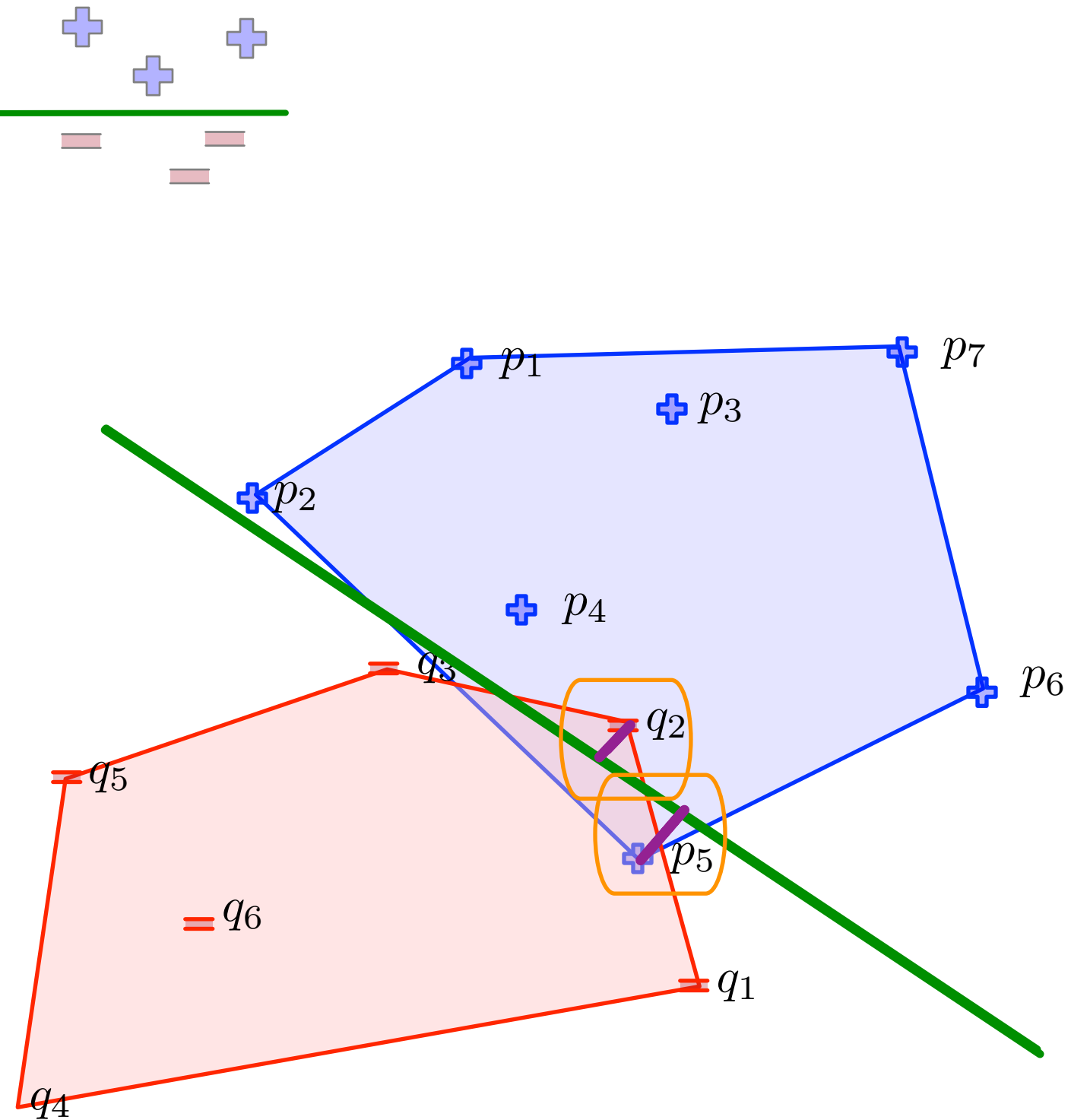
$\ell_i$

logistic

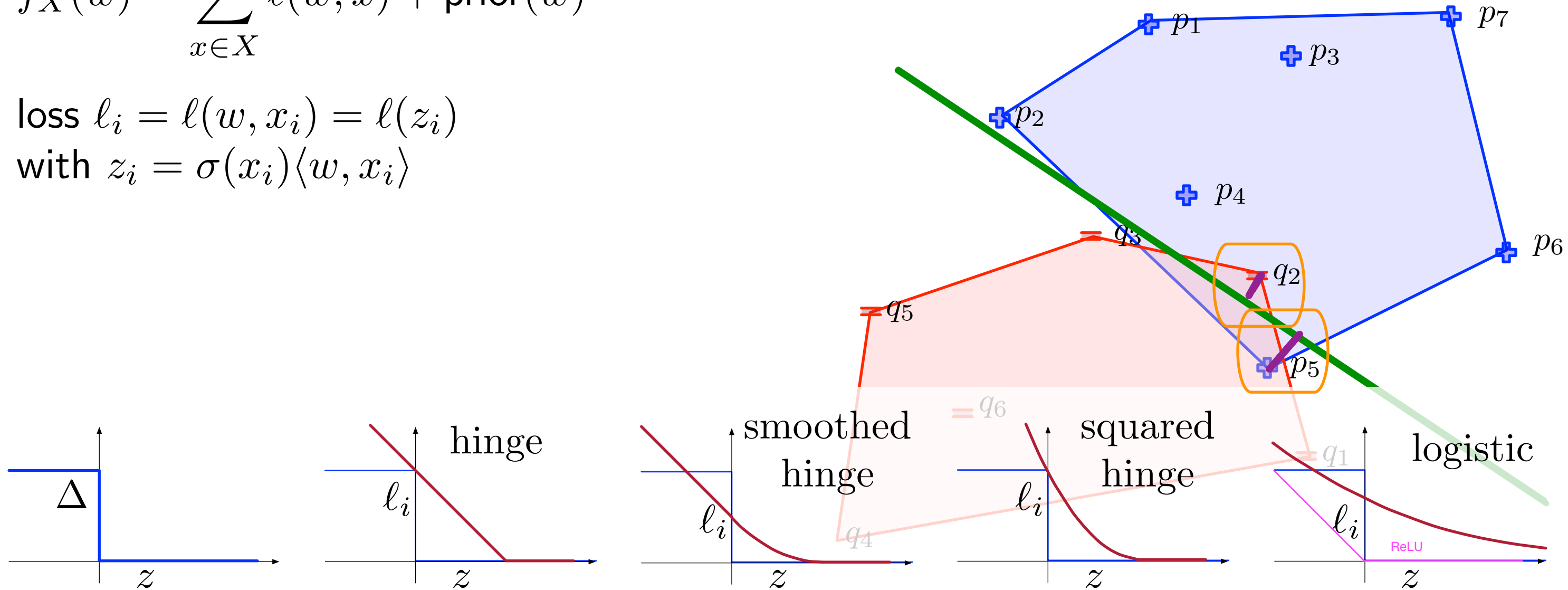$\ell_i$

ReLU

$z$  $z$  $z$  $z$  $z$

# Loss Functions

Set up penalty for misclassified points

$$f_X(w) = \sum_{x \in X} \ell(w, x) + \text{prior}(w)$$

loss $\ell_i = \ell(w, x_i) = \ell(z_i)$
with $z_i = \sigma(x_i)\langle w, x_i \rangle$

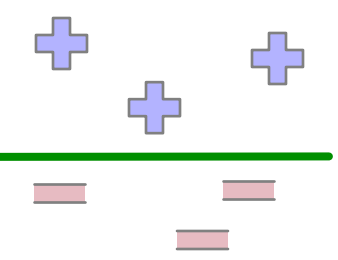If linear or $K$ is pd
then $f_X$ is **convex**

$p_1$  $p_7$
$p_3$
$p_2$
$p_4$
$q_3$
$q_2$
$q_5$
$p_5$
$q_6$
$q_4$
$q_1$

$\Delta$

hinge
$\ell_i$

smoothed
hinge
$\ell_i$

squared
hinge
$\ell_i$

logistic
$\ell_i$
ReLU

$z$   $z$   $z$   $z$   $z$

# Coresets for Optimization

Solve for $w^* = \arg\min_w f_X(w)$ ... gradient descent

gamma 0.30 | final grad 0.092

hinge

smoothed hinge

squared hinge

logistic

$\Delta$

$\ell_i$

$\ell_i$

$\ell_i$

$\ell_i$

ReLU

$z$

# Coresets for Optimization

Solve for $w^* = \arg\min_w f_X(w)$ ... gradient descent

Subgradient Descent $\approx$ Frank-Wolfe
Move towards most helpful point
$O(C/\varepsilon^2)$ steps to $\pm\varepsilon$ mean for $\|x\| \leq 1$.

Stochastic gradient descent (SGD)
For large $X$, most common
Randomly chooses $x \in X$, and step towards $-\nabla f_x(w)$.

gamma 0.30 | final grad 0.092

$\Delta$

$\ell_i$  hinge

$\ell_i$  smoothed hinge

$\ell_i$  squared hinge

$\ell_i$  logistic

ReLU

# VC-dimension and Sample Complexity

Assume: *data $X$ is drawn iid from $\mu$.*

Build: classifier $g : \mathbb{R}^d \to \{-1, +1\}$ so $\mathbf{E}_{x \sim \mu}[\mathbf{1}(g(x) = \sigma(x))]$.

# VC-dimension and Sample Complexity

Assume: *data $X$ is drawn iid from $\mu$.*

Build: classifier $g : \mathbb{R}^d \to \{-1, +1\}$ so $\mathbf{E}_{x \sim \mu}[\mathbf{1}(g(x) = \sigma(x))]$.

*how much accuracy is preserved under random sampling?*

# VC-dimension and Sample Complexity

Assume: *data $X$ is drawn iid from $\mu$.*

Build: classifier $g : \mathbb{R}^d \to \{-1, +1\}$ so $\mathbf{E}_{x \sim \mu}[\mathbf{1}(g(x) = \sigma(x))]$.

*how much accuracy is preserved under random sampling?*

**Range space** $(X, \mathcal{R})$, where $\mathcal{R}$ is family of subsets.

# VC-dimension and Sample Complexity

Assume: *data $X$ is drawn iid from $\mu$.*

Build: classifier $g : \mathbb{R}^d \rightarrow \{-1, +1\}$ so $\mathbf{E}_{x \sim \mu}[\mathbf{1}(g(x) = \sigma(x))]$.

*how much accuracy is preserved under random sampling?*

**Range space** $(X, \mathcal{R})$, where $\mathcal{R}$ is family of subsets.
**VC-dimension** $\nu$: largest $Y \subset X$ so all $Z \subset Y$ can be defined so $Z = Y \cap R$ for some $R \in \mathcal{R}$.

$\leftrightarrow$ sublinear range searching, approximate hitting sets

# VC-dimension and Sample Complexity

Assume: *data $X$ is drawn iid from $\mu$.*

Build: classifier $g : \mathbb{R}^d \to \{-1, +1\}$ so $\mathbf{E}_{x \sim \mu}[\mathbf{1}(g(x) = \sigma(x))]$.

*how much accuracy is preserved under random sampling?*

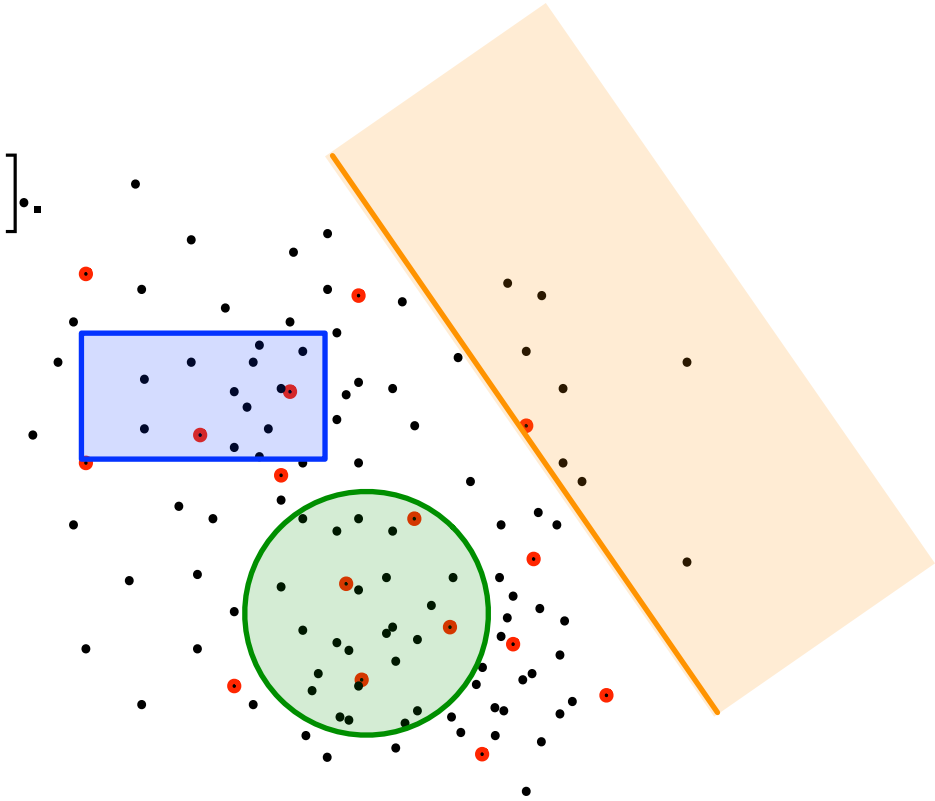**Range space** $(X, \mathcal{R})$, where $\mathcal{R}$ is family of subsets.
**VC-dimension** $\nu$: largest $Y \subset X$ so all $Z \subset Y$ can be defined so $Z = Y \cap R$ for some $R \in \mathcal{R}$.

$\leftrightarrow$ sublinear range searching, approximate hitting sets

An $\varepsilon$-**sample** is a subset $S \subset X$ so $|S| = O(1/\varepsilon^{2-2/(\nu+1)})$

$$\max_{R \in \mathcal{R}} \left| \frac{|X \cap R|}{|X|} - \frac{|S \cap R|}{|S|} \right| \leq \varepsilon.$$

# VC-dimension and Sample Complexity

Assume: *data $X$ is drawn iid from $\mu$.*
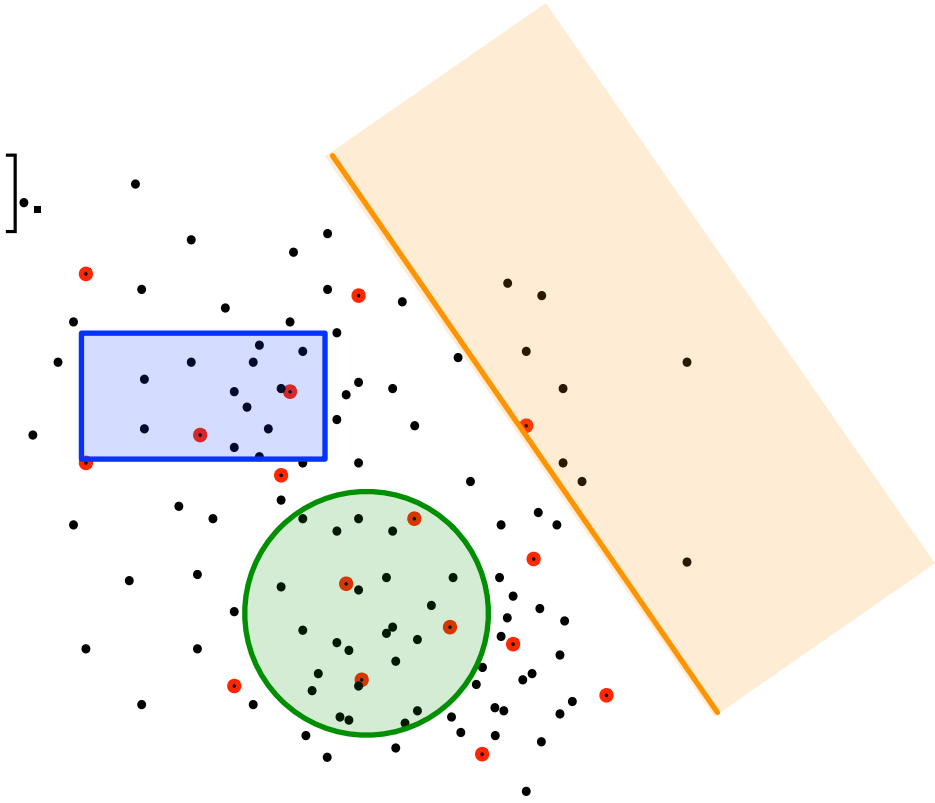
Build: classifier $g : \mathbb{R}^d \to \{-1, +1\}$ so $\mathbf{E}_{x \sim \mu}[\mathbf{1}(g(x) = \sigma(x))]$.
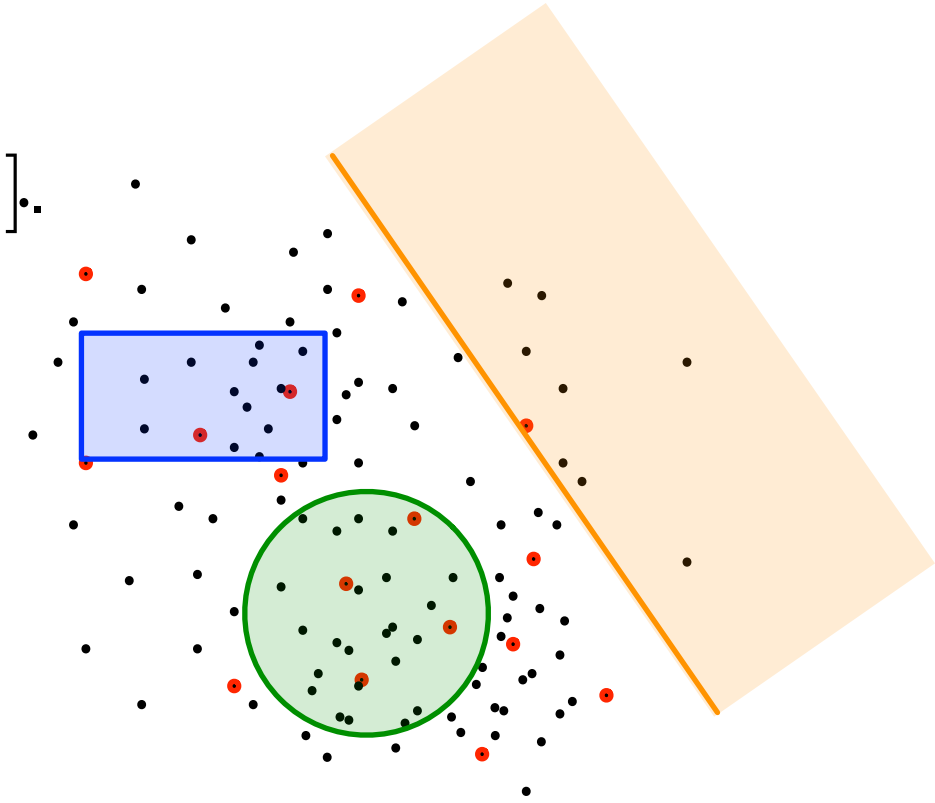
*how much accuracy is preserved under random sampling?*

**Range space** $(X, \mathcal{R})$, where $\mathcal{R}$ is family of subsets.
**VC-dimension** $\nu$: largest $Y \subset X$ so all $Z \subset Y$ can be
defined so $Z = Y \cap R$ for some $R \in \mathcal{R}$.

$\leftrightarrow$ sublinear range searching, approximate hitting sets

An $\varepsilon$-**sample** is a subset $S \subset X$ so
$$\max_{R \in \mathcal{R}} \left| \frac{|X \cap R|}{|X|} - \frac{|S \cap R|}{|S|} \right| \leq \varepsilon.$$

|  | **Geom** | **ML** |
|---|---|---|
| $\mathcal{R}$ | ranges | classifiers |
| $X$ | ground set | samples |

# VC-dimension and Sample Complexity

Assume: *data $X$ is drawn iid from $\mu$.*

Build: classifier $g : \mathbb{R}^d \to \{-1, +1\}$ so $\mathbf{E}_{x \sim \mu}[\mathbf{1}(g(x) = \sigma(x))]$.
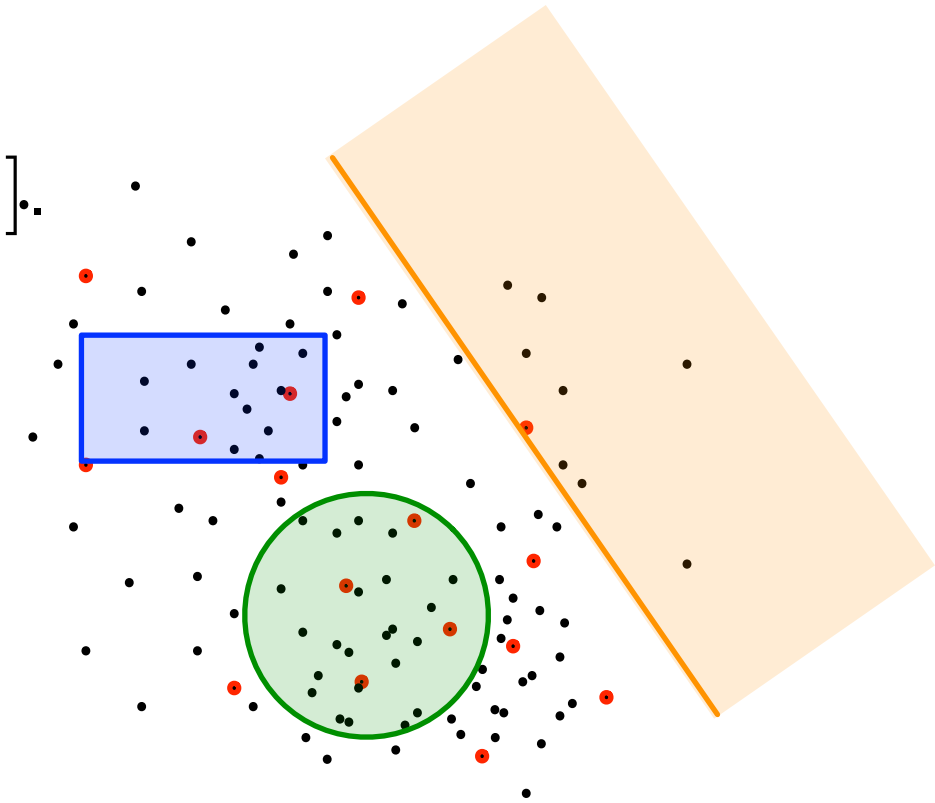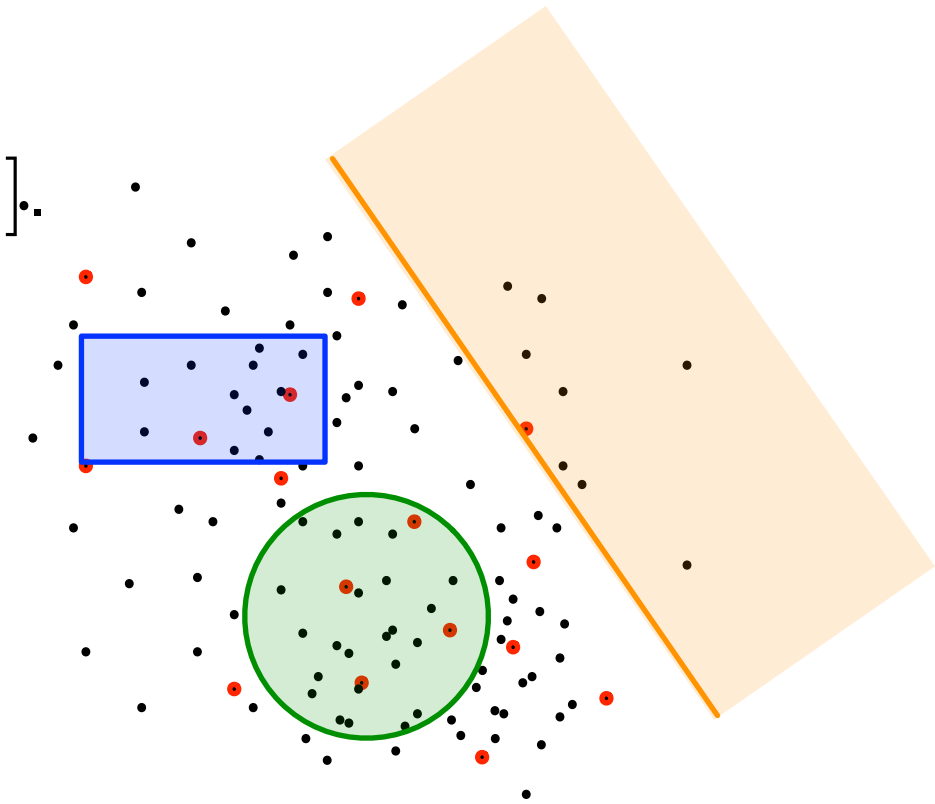
*how much accuracy is preserved under random sampling?*

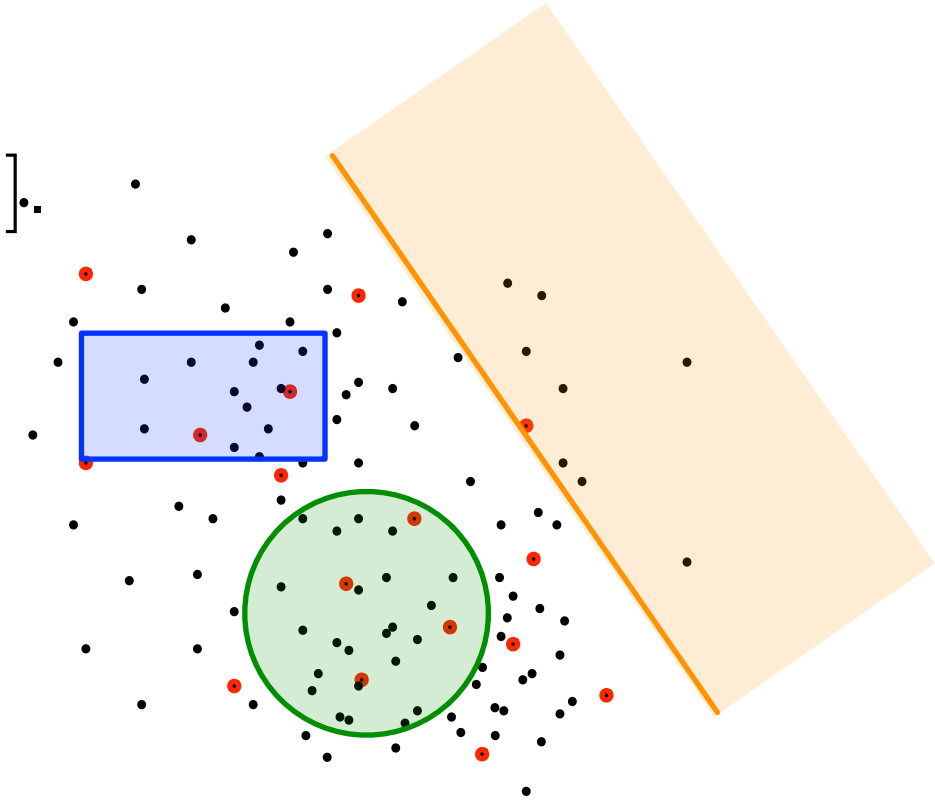**Range space** $(X, \mathcal{R})$, where $\mathcal{R}$ is family of subsets.
**VC-dimension** $\nu$: largest $Y \subset X$ so all $Z \subset Y$ can be
defined so $Z = Y \cap R$ for some $R \in \mathcal{R}$.

$\leftrightarrow$ sublinear range searching, approximate hitting sets

An $\varepsilon$-**sample** is a subset $S \subset X$ so
$$\max_{R \in \mathcal{R}} \left| \frac{|X \cap R|}{|X|} - \frac{|S \cap R|}{|S|} \right| \leq \varepsilon.$$

| | **Geom** | **ML** |
|---|---|---|
| $\mathcal{R}$ | ranges | classifiers |
| $X$ | ground set | samples |
| $O(\frac{\nu}{\varepsilon^2})$ | $\varepsilon$-sample [VC71] | agnostic learning |

# VC-dimension and Sample Complexity

Assume: *data $X$ is drawn iid from $\mu$.*

Build: classifier $g : \mathbb{R}^d \to \{-1, +1\}$ so $\mathbf{E}_{x \sim \mu}[\mathbf{1}(g(x) = \sigma(x))]$.

*how much accuracy is preserved under random sampling?*

**Range space** $(X, \mathcal{R})$, where $\mathcal{R}$ is family of subsets.
**VC-dimension** $\nu$: largest $Y \subset X$ so all $Z \subset Y$ can be defined so $Z = Y \cap R$ for some $R \in \mathcal{R}$.

$\leftrightarrow$ sublinear range searching, approximate hitting sets

An $\varepsilon$-**sample** is a subset $S \subset X$ so

$$\max_{R \in \mathcal{R}} \left| \frac{|X \cap R|}{|X|} - \frac{|S \cap R|}{|S|} \right| \leq \varepsilon.$$

An $\varepsilon$-**net** is a subset $S \subset X$ such that each $R \in \mathcal{R}$ so $\frac{|R \cap X|}{|X|} \geq \varepsilon$ then $S \cap R \neq \emptyset$.
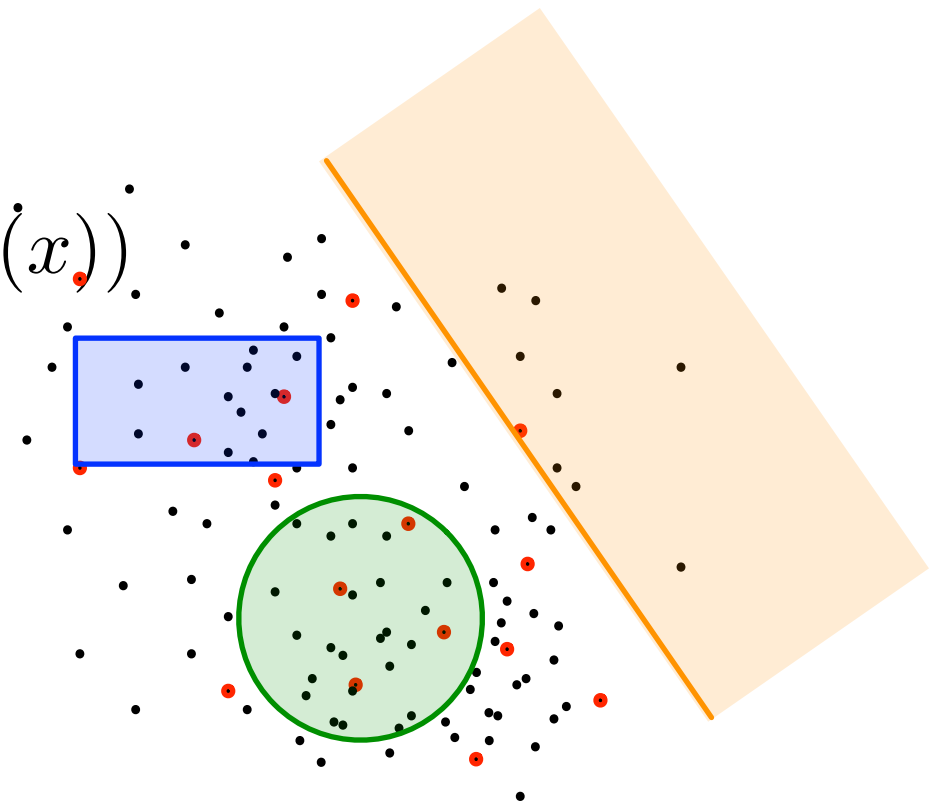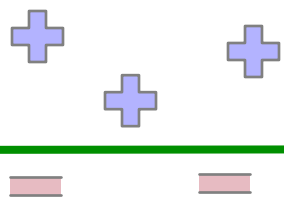
|  | **Geom** | **ML** |
|---|---|---|
| $\mathcal{R}$ | ranges | classifiers |
| $X$ | ground set | samples |
| $O(\frac{\nu}{\varepsilon^2})$ | $\varepsilon$-sample [VC71] | agnostic learning |
| $O(\frac{\nu}{\varepsilon} \log \frac{\nu}{\varepsilon})$ | $\varepsilon$-net [HW85] | perfect classifiers |

# **Recent Results** (... of mine)

- Approximate Maximum Disagreement Problem $X, \overline{\sigma}, \mathcal{H}$

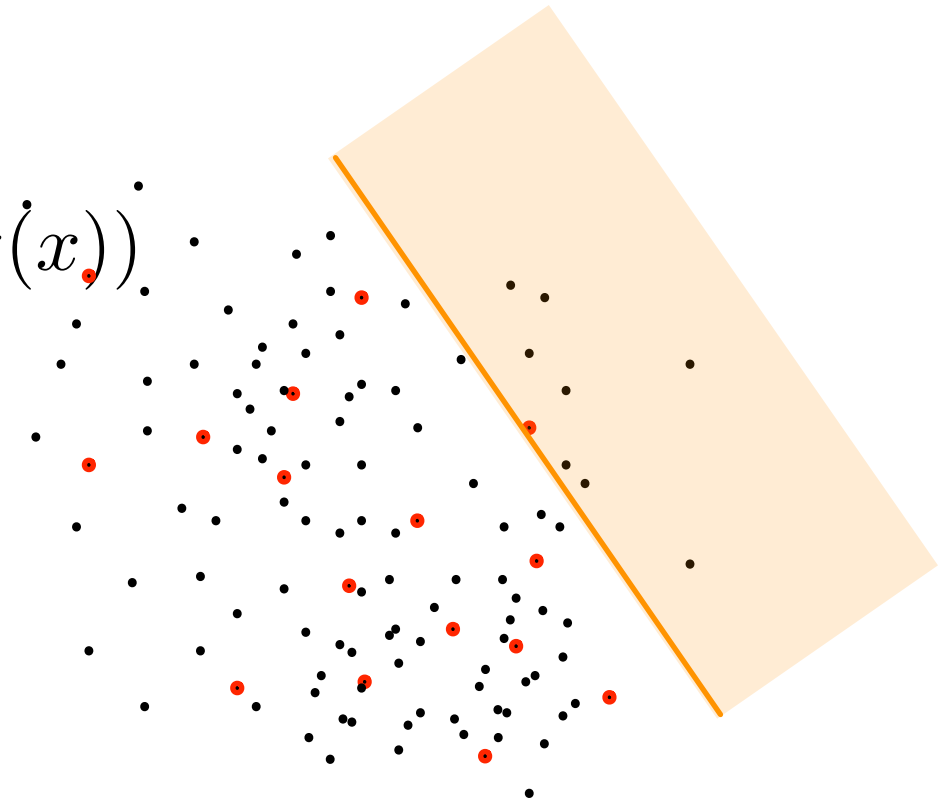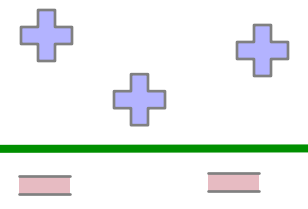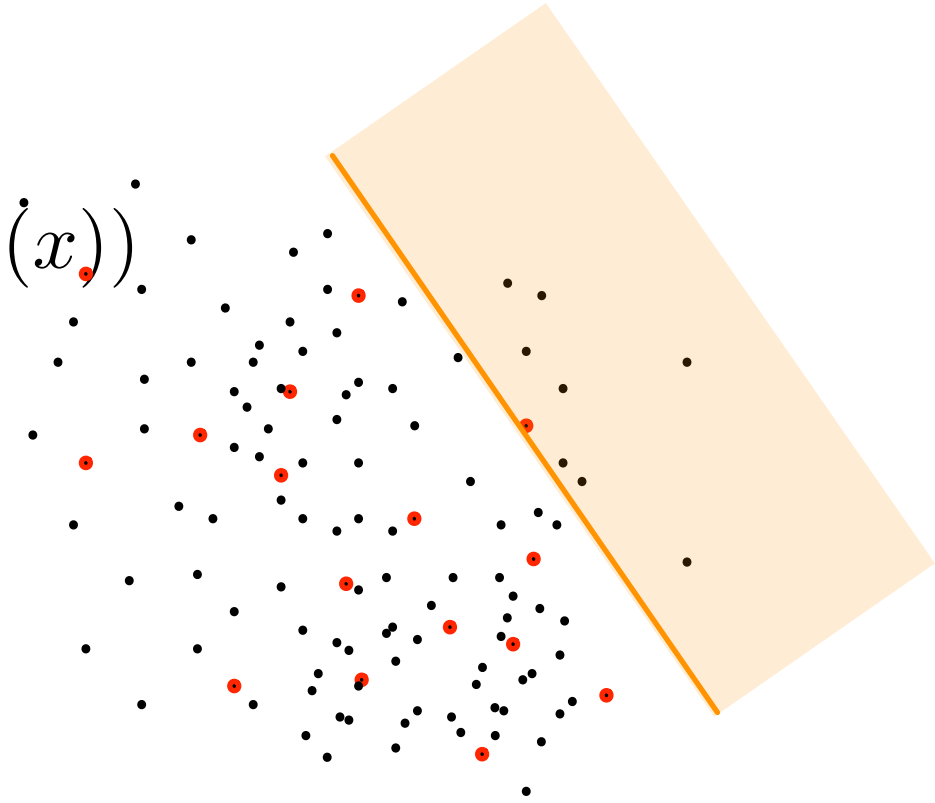Find halfspace $h$ maximizes $\Delta_X(h) = \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(h(x) = \sigma(x))$

# Recent Results (... of mine)

- Approximate Maximum Disagreement Problem $X, \overline{\sigma}, \mathcal{H}$

Find halfspace $h$ maximizes $\Delta_X(h) = \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(h(x) = \sigma(x))$

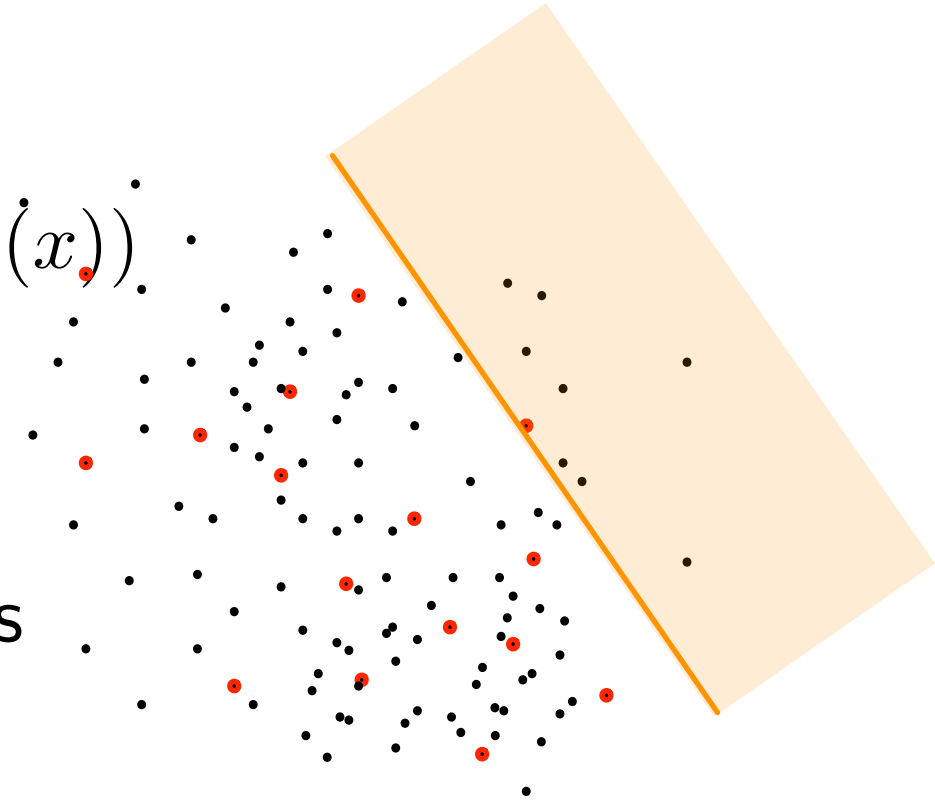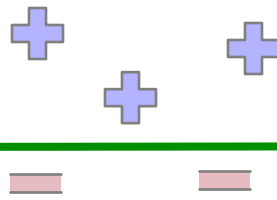Find $\hat{h} \in \mathcal{H}$ so $\Delta_X(h^*) - \Delta_X(\hat{h}) \leq \varepsilon$

# Recent Results (... of mine)

• Approximate Maximum Disagreement Problem $X, \bar{\sigma}, \mathcal{H}$

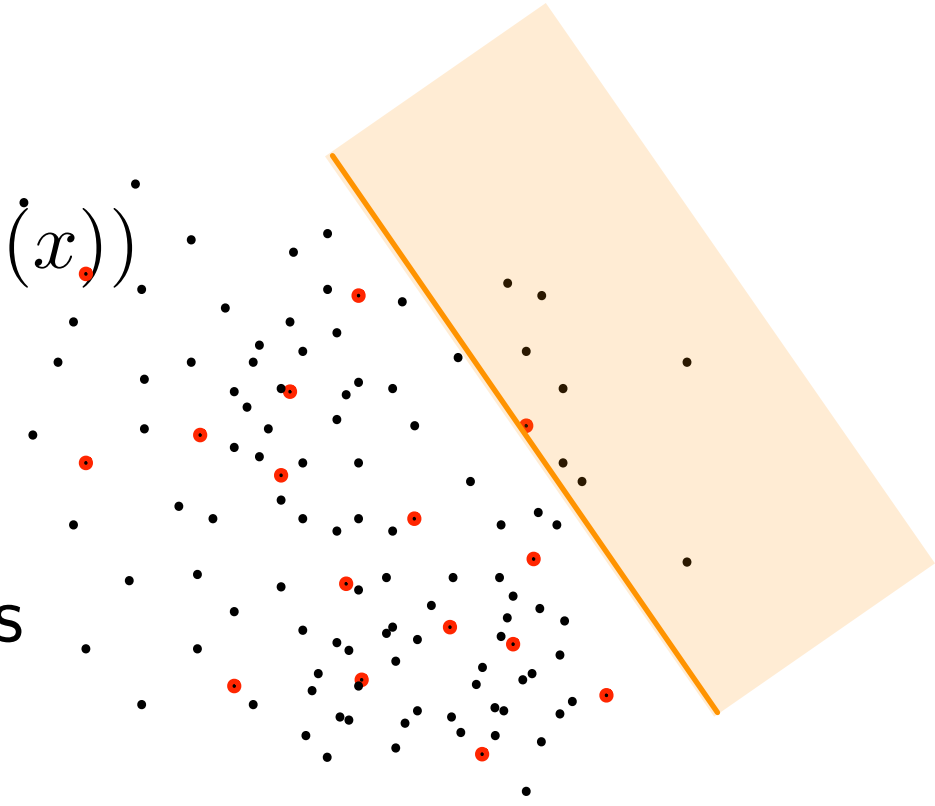Find halfspace $h$ maximizes $\Delta_X(h) = \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(h(x) = \sigma(x))$

Find $\hat{h} \in \mathcal{H}$ so $\Delta_X(h^*) - \Delta_X(\hat{h}) \leq \varepsilon$

... in $O((1/\varepsilon^{d+1/3})\text{polylog}(1/\varepsilon))$ time. [Matheny+P '18]

# Recent Results (... of mine)

- Approximate Maximum Disagreement Problem $X, \bar{\sigma}, \mathcal{H}$

Find halfspace $h$ maximizes $\Delta_X(h) = \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(h(x) = \sigma(x))$

Find $\hat{h} \in \mathcal{H}$ so $\Delta_X(h^*) - \Delta_X(\hat{h}) \leq \varepsilon$

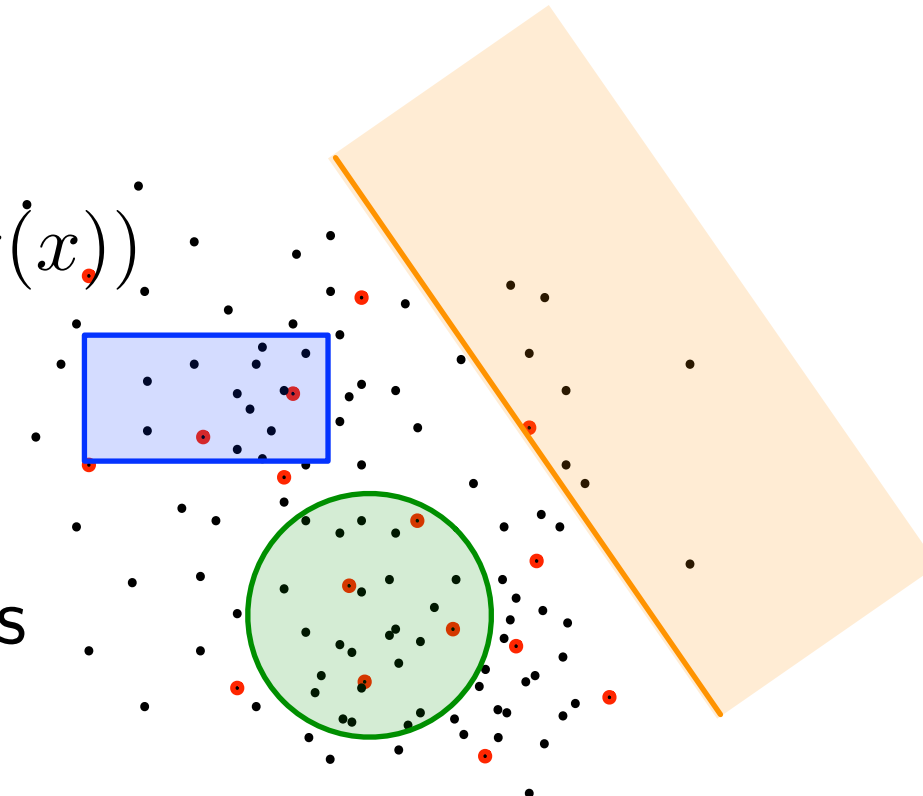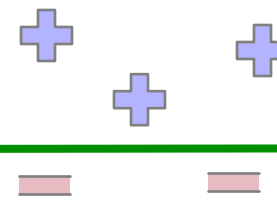... in $O((1/\varepsilon^{d+1/3})\text{polylog}(1/\varepsilon))$ time. [Matheny+P '18]

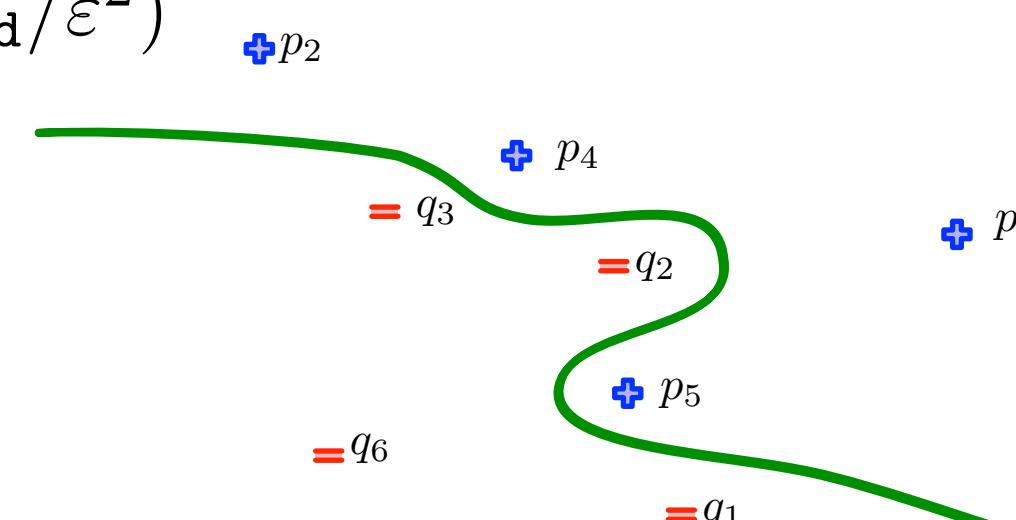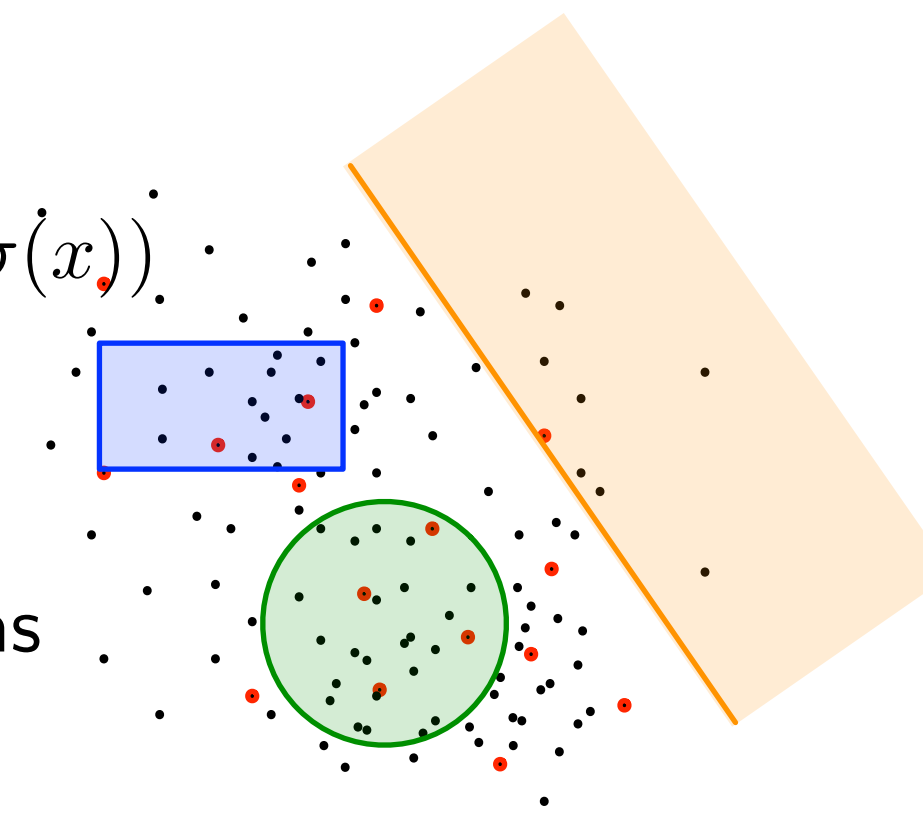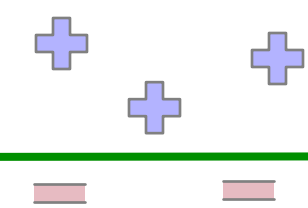- Sublevel-Set Range Space $(X, \mathcal{R}_d)$ defined for $R_{p,\tau} \in \mathcal{R}_d$ as

$R_{p,\tau} = \{x \in X \mid \mathsf{d}(p, x) \leq \tau\}$.

# Recent Results (... of mine)

- Approximate Maximum Disagreement Problem $X, \overline{\sigma}, \mathcal{H}$

Find halfspace $h$ maximizes $\Delta_X(h) = \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(h(x) = \sigma(x))$

Find $\hat{h} \in \mathcal{H}$ so $\Delta_X(h^*) - \Delta_X(\hat{h}) \leq \varepsilon$

... in $O((1/\varepsilon^{d+1/3})\text{polylog}(1/\varepsilon))$ time. [Matheny+P '18]

- Sublevel-Set Range Space $(X, \mathcal{R}_\mathsf{d})$ defined for $R_{p,\tau} \in \mathcal{R}_\mathsf{d}$ as

$R_{p,\tau} = \{x \in X \mid \mathsf{d}(p,x) \leq \tau\}$.

General KDE: $K_\mathsf{d}(x,p) = \exp(-\mathsf{d}(x,p)^2)$ and $\text{KDE}_X(p) = \frac{1}{|X|} \sum_{x \in X} K(x,p)$

# Recent Results (... of mine)

- Approximate Maximum Disagreement Problem $X, \bar{\sigma}, \mathcal{H}$

Find halfspace $h$ maximizes $\Delta_X(h) = \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(h(x) = \sigma(x))$

Find $\hat{h} \in \mathcal{H}$ so $\Delta_X(h^*) - \Delta_X(\hat{h}) \leq \varepsilon$

... in $O((1/\varepsilon^{d+1/3})\text{polylog}(1/\varepsilon))$ time. [Matheny+P '18]

- Sublevel-Set Range Space $(X, \mathcal{R}_\mathsf{d})$ defined for $R_{p,\tau} \in \mathcal{R}_\mathsf{d}$ as
$R_{p,\tau} = \{x \in X \mid \mathsf{d}(p,x) \leq \tau\}$.

General KDE: $K_\mathsf{d}(x,p) = \exp(-\mathsf{d}(x,p)^2)$ and $\text{KDE}_X(p) = \frac{1}{|X|} \sum_{x \in X} K(x,p)$

If VC-dimension of $(X, \mathcal{R}_\mathsf{d})$ is $\nu_\mathsf{d}$ then $S \sim X$ of size $O(\nu_\mathsf{d}/\varepsilon^2)$
has $\|\text{KDE}_X - \text{KDE}_S\|_\infty \leq \varepsilon$ [Komaraju etal '11]

# Recent Results (... of mine)

• Approximate Maximum Disagreement Problem $X, \sigma, \mathcal{H}$
Find halfspace $h$ maximizes $\Delta_X(h) = \frac{1}{|X|} \sum_{x \in X} \mathbf{1}(h(x) = \sigma(x))$

Find $\hat{h} \in \mathcal{H}$ so $\Delta_X(h^*) - \Delta_X(\hat{h}) \leq \varepsilon$

... in $O((1/\varepsilon^{d+1/3})\text{polylog}(1/\varepsilon))$ time. [Matheny+P '18]

• Sublevel-Set Range Space $(X, \mathcal{R}_\mathsf{d})$ defined for $R_{p,\tau} \in \mathcal{R}_\mathsf{d}$ as
$R_{p,\tau} = \{x \in X \mid \mathsf{d}(p,x) \leq \tau\}$.

General KDE: $K_\mathsf{d}(x,p) = \exp(-\mathsf{d}(x,p)^2)$ and $\text{KDE}_X(p) = \frac{1}{|X|} \sum_{x \in X} K(x,p)$

If VC-dimension of $(X, \mathcal{R}_\mathsf{d})$ is $\nu_\mathsf{d}$ then $S \sim X$ of size $O(\nu_\mathsf{d}/\varepsilon^2)$
has $\|\text{KDE}_X - \text{KDE}_S\|_\infty \leq \varepsilon$ [Komaraju etal '11]
approximates margin $\pm\varepsilon$.

# Many, many types of classifiers

- **KNN**: voting based on $k$ nearest neighbors

# Many, many types of classifiers

- **KNN**: voting based on $k$ nearest neighbors

- **decision trees**: kd-tree until uniform nodes

$p_1$

$p_7$

$p_3$

$p_2$

$p_4$

$q_3$

$p_6$

$q_2$

$q_5$

$p_5$

$q_6$

$q_1$

$q_4$

# Many, many types of classifiers

- **KNN**: voting based on $k$ nearest neighbors

- **decision trees**: kd-tree until uniform nodes

- **random forests**: mixtures of decision trees

$p_1$

$p_7$

$p_3$

$p_2$

$p_4$

$q_3$

$p_6$

$q_2$

$q_5$

$p_5$

$q_6$

$q_1$

$q_4$

# Many, many types of classifiers

- **KNN**: voting based on $k$ nearest neighbors

- **decision trees**: kd-tree until uniform nodes

- **random forests**: mixtures of decision trees

- **deep neural networks**: many layers of perceptrons
piecewise-linear classifiers

input layer

hidden layer 1   hidden layer 2   hidden layer 3

output layer

$p_1$   $p_7$

$p_3$

$p_2$

$p_4$

$q_3$   $p_6$

$q_2$

$q_5$

$p_5$

$q_6$

$q_1$

$q_4$

# Many, many types of classifiers

- **KNN**: voting based on $k$ nearest neighbors

- **decision trees**: kd-tree until uniform nodes

- **random forests**: mixtures of decision trees

- **deep neural networks**: many layers of perceptrons piecewise-linear classifiers

- non-convex function $f_X(w)$
- lots of data!
- lots of computers!

$p_3$

input layer   hidden layer 1   hidden layer 2   hidden layer 3

output layer

$= q_5$

$= q_4$

# Many, many types of classifiers

- **KNN**: voting based on $k$ nearest neighbors

- **decision trees**: kd-tree until uniform nodes

- **random forests**: mixtures of decision trees

- **deep neural networks**: many layers of perceptrons
piecewise-linear classifiers

- ...

$p_1$  $p_7$

$p_3$

$p_2$

$p_4$

$q_3$  $p_6$

$q_2$

$q_5$

$p_5$

$q_6$

$q_1$

$q_4$

# Where can geometry help today?

- non-convex optimization, function classes

- generalizations of VC-dimension

- (implicit) regularization:
what it means, when it occurs

- geometry of data ... e.g., adversarial data

- bias in data, algorithms

# Where can geometry help today?

- non-convex optimization, function classes

- generalizations of VC-dimension

- (implicit) regularization:
what it means, when it occurs

- geometry of data … e.g., adversarial data

- bias in data, algorithms

**word vector embeddings**
word2vec, GloVe, …

doctor

man

king

woman

queen

nurse

car

truck

$\mathbb{R}^{300}$

# Where can geometry help today?

- non-convex optimization, function classes

- generalizations of VC-dimension

- (implicit) regularization:
what it means, when it occurs

- geometry of data ... e.g., adversarial data

- bias in data, algorithms

**word vector embeddings**
word2vec, GloVe, ...

# Where can geometry help today?

- non-convex optimization, function classes

- generalizations of VC-dimension

- (implicit) regularization:
what it means, when it occurs

- geometry of data ... e.g., adversarial data

- bias in data, algorithms

**word vector embeddings**
word2vec, GloVe, ...

doctor

man

king

woman

queen

nurse

car

truck

$\mathbb{R}^{300}$

# Where can geometry help today?

- non-convex optimization, function classes

- generalizations of VC-dimension

- (implicit) regularization:
what it means, when it occurs

- geometry of data ... e.g., adversarial data

- bias in data, algorithms

**word vector embeddings**
word2vec, GloVe, ...

doctor

man

king

woman

*gender*

queen

nurse

car

truck

$\mathbb{R}^{300}$

# Where can geometry help today?

- non-convex optimization, function classes

- generalizations of VC-dimension

- (implicit) regularization:
what it means, when it occurs

- geometry of data ... e.g., adversarial data

- bias in data, algorithms

**word vector embeddings**
word2vec, GloVe, ...

# Where can geometry help today?

- non-convex optimization, function classes

- generalizations of VC-dimension

- (implicit) regularization:
what it means, when it occurs

- geometry of data ... e.g., adversarial data

- bias in data, algorithms

**word vector embeddings**
word2vec, GloVe, ...

doctor

man

woman

*gender*

nurse

king

queen

car

truck

nurse
doctor

[Dev+P AIStats19]

$\mathbb{R}^{299}$

$\mathbb{R}^{300}$

# Mathematical Foundations for Data Analysis

## Jeff M. Phillips

http://www.cs.utah.edu/~jeffp/M4D/M4D.html