

Learning a set of directions



Wouter M. Koolen Jiazhong Nie Manfred K. Warmuth

Outline

Motivation

Measuring gain

Algorithm

Conclusion

Problem

Parts of my home town Amsterdam lie 5 metres below sea level



Solution

Pump out water



Leeghwater (1607)

This is how we do it



And then global warming sets in ...



Online learning to the rescue

For $t = 1, 2, \dots$

- ▶ Mill chooses a direction u_t
- ▶ Wind reveals direction x_t
- ▶ Gain based on match

What is a reasonable gain?

Measuring gain

Gain quantifies quality of prediction u upon outcome x

Measuring gain

Gain quantifies quality of prediction u upon outcome x

Perhaps the simplest gain is the

$$\text{angle cosine} := u^T x$$

Measuring gain

Gain quantifies quality of prediction u upon outcome x

Perhaps the simplest gain is the

$$\text{angle cosine} := u^T x$$

best when u, x parallel, worst when u, x opposite

Measuring gain

Gain quantifies quality of prediction u upon outcome x

Perhaps the simplest gain is the

$$\text{angle cosine} := u^T x$$

best when u, x parallel, worst when u, x opposite

Another gain is used in Principal Component Analysis

$$\text{subspace similarity} := (u^T x)^2$$

Measuring gain

Gain quantifies quality of prediction u upon outcome x

Perhaps the simplest gain is the

$$\text{angle cosine} := u^T x$$

best when u, x parallel, worst when u, x opposite

Another gain is used in Principal Component Analysis

$$\text{subspace similarity} := (u^T x)^2$$

best when u, x parallel or opposite

Measuring gain

Gain quantifies quality of prediction u upon outcome x

Perhaps the simplest gain is the

$$\text{angle cosine} := u^T x$$

best when u, x parallel, **worst when u, x opposite**

Another gain is used in Principal Component Analysis

$$\text{subspace similarity} := (u^T x)^2$$

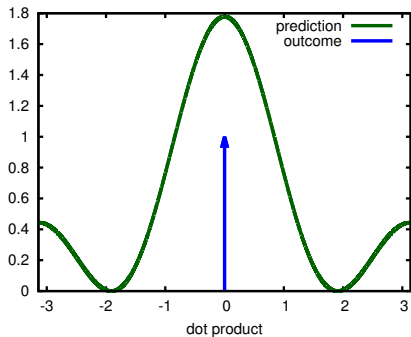
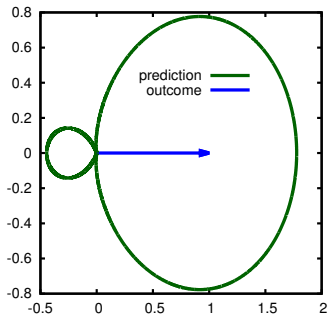
best when u, x parallel **or opposite**

Our solution: controlled trade-off (windmill-dependent constant c)

$$\text{directional gain} := (u^T x + c)^2$$

Visualisation of directional gain

$$(\mathbf{u}^T \mathbf{x} + c)^2 \quad \text{with } c = 1/3$$



Gain expansion

For randomised prediction $\mathbf{u} \sim \mathbb{P}$:

$$\begin{aligned}\mathbb{E} \left[(\mathbf{u}^\top \mathbf{x} + c)^2 \right] &= \mathbb{E} \left[\mathbf{x}^\top \mathbf{u} \mathbf{u}^\top \mathbf{x} + 2c \mathbf{x}^\top \mathbf{u} + c^2 \right] \\ &= \mathbf{x}^\top \mathbb{E} [\mathbf{u} \mathbf{u}^\top] \mathbf{x} + 2c \mathbf{x}^\top \mathbb{E} [\mathbf{u}] + c^2.\end{aligned}$$

Only relevant characteristics of \mathbb{P} are its

$$\boldsymbol{\mu} := \mathbb{E} [\mathbf{u}]$$

first moment vector

$$\mathbf{D} := \mathbb{E} [\mathbf{u} \mathbf{u}^\top]$$

second moment matrix

Observation: gain is **linear** in $\boldsymbol{\mu}$ and in \mathbf{D}

Gain expansion

For randomised prediction $\mathbf{u} \sim \mathbb{P}$:

$$\begin{aligned}\mathbb{E} \left[(\mathbf{u}^\top \mathbf{x} + c)^2 \right] &= \mathbb{E} \left[\mathbf{x}^\top \mathbf{u} \mathbf{u}^\top \mathbf{x} + 2c \mathbf{x}^\top \mathbf{u} + c^2 \right] \\ &= \mathbf{x}^\top \mathbb{E} [\mathbf{u} \mathbf{u}^\top] \mathbf{x} + 2c \mathbf{x}^\top \mathbb{E} [\mathbf{u}] + c^2.\end{aligned}$$

Only relevant characteristics of \mathbb{P} are its

$$\boldsymbol{\mu} := \mathbb{E} [\mathbf{u}]$$

first moment vector

$$\mathbf{D} := \mathbb{E} [\mathbf{u} \mathbf{u}^\top]$$

second moment matrix

Observation: gain is **linear** in $\boldsymbol{\mu}$ and in \mathbf{D}

Idea: forget \mathbb{P} - use $\boldsymbol{\mu}$ and \mathbf{D} as a parameter

Gain expansion

For randomised prediction $\mathbf{u} \sim \mathbb{P}$:

$$\begin{aligned}\mathbb{E} \left[(\mathbf{u}^\top \mathbf{x} + c)^2 \right] &= \mathbb{E} \left[\mathbf{x}^\top \mathbf{u} \mathbf{u}^\top \mathbf{x} + 2c \mathbf{x}^\top \mathbf{u} + c^2 \right] \\ &= \mathbf{x}^\top \mathbb{E} [\mathbf{u} \mathbf{u}^\top] \mathbf{x} + 2c \mathbf{x}^\top \mathbb{E} [\mathbf{u}] + c^2.\end{aligned}$$

Only relevant characteristics of \mathbb{P} are its

$$\begin{array}{ll}\boldsymbol{\mu} := \mathbb{E} [\mathbf{u}] & \text{first moment vector} \\ \mathbf{D} := \mathbb{E} [\mathbf{u} \mathbf{u}^\top] & \text{second moment matrix}\end{array}$$

Observation: gain is **linear** in $\boldsymbol{\mu}$ and in \mathbf{D}

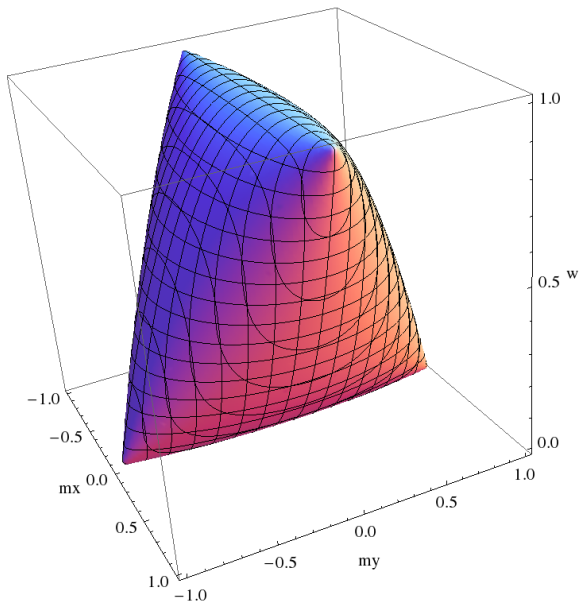
Idea: forget \mathbb{P} - use $\boldsymbol{\mu}$ and \mathbf{D} as a parameter

Careful: not all $\langle \boldsymbol{\mu}, \mathbf{D} \rangle$ are moments of some \mathbb{P} !

Parameters must lie in the **überplex**

$$\mathcal{U} := \{ \langle \boldsymbol{\mu}, \mathbf{D} \rangle \mid \exists \mathbb{P} : \boldsymbol{\mu}, \mathbf{D} \text{ are 1}^{\text{st}}/2^{\text{nd}} \text{ moment of } \mathbb{P} \}$$

Überplex \mathcal{U}



Characterisation

Theorem

$$\langle \boldsymbol{\mu}, \mathbf{D} \rangle \in \mathcal{U} \quad \text{iff} \quad \text{tr}(\mathbf{D}) = 1 \quad \text{and} \quad \mathbf{D} \succeq \boldsymbol{\mu}\boldsymbol{\mu}^\top$$

Characterisation

Theorem

$$\langle \boldsymbol{\mu}, \mathbf{D} \rangle \in \mathcal{U} \quad \text{iff} \quad \text{tr}(\mathbf{D}) = 1 \quad \text{and} \quad \mathbf{D} \succeq \boldsymbol{\mu}\boldsymbol{\mu}^\top$$

Why this is important?

- ▶ Überplex \mathcal{U} is convex
- ▶ Constraint is semi-definite
- ▶ Efficient numerical linear/convex optimization over \mathcal{U}

Offline problem

$$\max_{(\boldsymbol{\mu}, \mathbf{D}) \in \mathcal{U}} \sum_{t=1}^T (\mathbf{x}_t^T \mathbf{D} \mathbf{x}_t + 2c\boldsymbol{\mu}^T \mathbf{x}_t + c^2)$$

Semi-definite optimisation problem
Good numerical methods

“Our” algorithm: gradient descent

Maintains the two moments $(\boldsymbol{\mu}_t, \mathbf{D}_t) \in \mathcal{U}$ as parameter

At trial $t = 1 \dots T$

1. Mill decomposes parameter $(\boldsymbol{\mu}_t, \mathbf{D}_t)$ into a mixture of 6 directions and draws a direction \mathbf{u}_t at random from it
2. Wind reveals direction $\mathbf{x}_t \in \mathbb{R}^2$
3. Mill receives expected gain $\mathbb{E} [(\mathbf{u}_t^\top \mathbf{x}_t + c)^2]$
4. Mill updates $(\boldsymbol{\mu}_t, \mathbf{D}_t)$ to $(\hat{\boldsymbol{\mu}}_{t+1}, \hat{\mathbf{D}}_{t+1})$ based on the gradient of the expected gain on \mathbf{x}_t

$$\hat{\boldsymbol{\mu}}_{t+1} := \boldsymbol{\mu}_t + 2\eta c \mathbf{x}_t \quad \text{and} \quad \hat{\mathbf{D}}_{t+1} := \mathbf{D}_t + \eta \mathbf{x}_t \mathbf{x}_t^\top$$

5. Mill produces new parameter $(\boldsymbol{\mu}_{t+1}, \mathbf{D}_{t+1})$ by projecting $(\hat{\boldsymbol{\mu}}_{t+1}, \hat{\mathbf{D}}_{t+1})$ back into the überplex

$$(\boldsymbol{\mu}_{t+1}, \mathbf{D}_{t+1}) := \underset{(\boldsymbol{\mu}, \mathbf{D}) \in \mathcal{U}}{\operatorname{argmin}} \|\mathbf{D} - \hat{\mathbf{D}}_{t+1}\|_F^2 + \|\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}_{t+1}\|^2$$

Guarantees

regret := hindsight-optimal gain – actual gain of Mill

Theorem

The expected regret after T trials of the GD algorithm with learning rate $\eta = \sqrt{\frac{3/2}{(4c^2+1)T}}$ and initial parameters $\mu_1 = \mathbf{0}$ and $D_1 = \frac{1}{2}\mathbf{I}$ is upper bounded by $\sqrt{3(4c^2 + 1)T}$

Guarantees

regret := hindsight-optimal gain – actual gain of **Mill**

Theorem

The expected regret after T trials of the GD algorithm with learning rate $\eta = \sqrt{\frac{3/2}{(4c^2+1)T}}$ and initial parameters $\boldsymbol{\mu}_1 = \mathbf{0}$ and $\mathbf{D}_1 = \frac{1}{2}\mathbf{I}$ is upper bounded by $\sqrt{3(4c^2 + 1)T}$

- ▶ Regret grows sub-linearly with T
- ▶ **Mill** turned close to the best orientation
- ▶ Holland is saved 😊

Conclusion

- ▶ An efficient method for orienting windmills
- ▶ Characterization of set of first two moments of distributions on directions

We can do more

- ▶ Work in $n \geq 3$ dimensions
- ▶ Learn sets of $k \geq 1$ orthogonal directions

